

# EXHIBIT F

**Family list**

**2** family member for:

**JP6083386**

Derived from 1 application.

**1 SPEECH RECOGNITION DEVICE FOR UNSPECIFIED  
SPEAKER**

Publication info: **JP3285048B2 B2** - 2002-05-27

**JP6083386 A** - 1994-03-25

---

Data supplied from the **esp@cenet** database - Worldwide

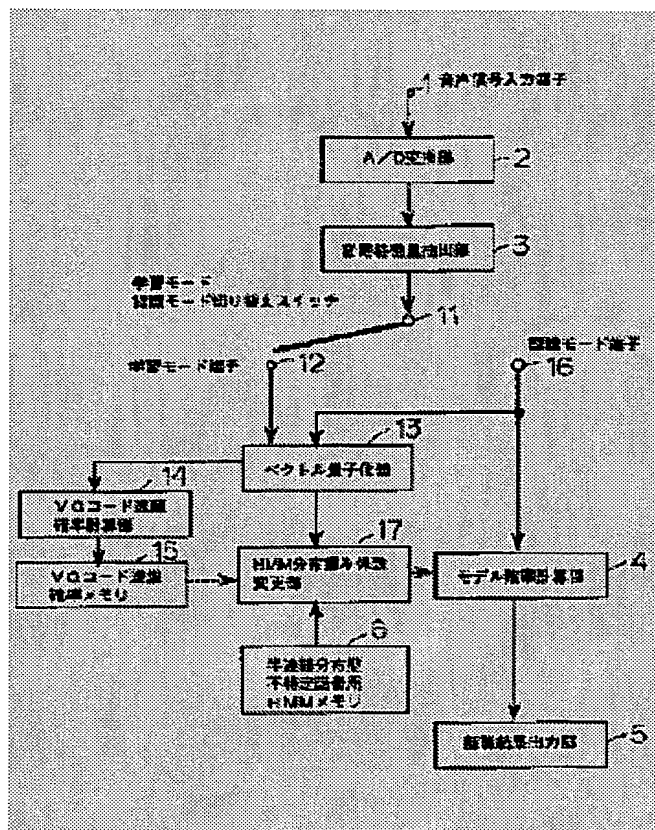
# SPEECH RECOGNITION DEVICE FOR UNSPECIFIED SPEAKER

**Patent number:** JP6083386  
**Publication date:** 1994-03-25  
**Inventor:** TAKAHASHI SATOSHI; others: 02  
**Applicant:** NIPPON TELEGR & TELEPH CORP  
**Classification:**  
 - international: G10L3/00  
 - european:  
**Application number:** JP19920237153 19920904  
**Priority number(s):**

## Abstract of JP6083386

**PURPOSE:** To reduce an overlap of distributions of models and improve the recognition performance by varying the weight of the distributions so that it matches an input speaker and phonemes by using the connection probability of vector quantization codes between two frames.

**CONSTITUTION:** A learning mode/recognition mode changeover switch 11 is set to the side of a recognition mode terminal 16 and the input speaker vocalizes a word to be recognized. The input speech is converted by a vector quantizer 13 into a series of codes, connection probability  $p(c_j|c_i)$  corresponding to the series of codes is read out of a vector quantized code connection probability memory 15, and an HMM distribution weight coefficient variation part 17 calculates the distribution weight coefficient  $\lambda$  of a half-connection distribution type unspecified speaker HMM (Hidden Markov Model) read out of a memory 6 by using the  $p(c_j|c_i)$  and varies it into  $\lambda_{daj}$ . Then calculation is performed by using the coefficient  $\lambda_{dajm}$  and a model probability calculation part 4 calculates probability for each HMM of the input speech on the basis of the calculation and other HMM parameters.



(51)Int.Cl.<sup>5</sup>

G 1 0 L 3/00

識別記号

5 3 5

5 3 1 J

庁内整理番号

7627-5H

7627-5H

F I

技術表示箇所

審査請求 未請求 請求項の数1(全5頁)

(21)出願番号

特願平4-237153

(22)出願日

平成4年(1992)9月4日

(71)出願人 000004226

日本電信電話株式会社

東京都千代田区内幸町一丁目1番6号

(72)発明者 高橋 敏

東京都千代田区内幸町1丁目1番6号 日

本電信電話株式会社内

(72)発明者 鹿野 清宏

東京都千代田区内幸町1丁目1番6号 日

本電信電話株式会社内

(72)発明者 松岡 達雄

東京都千代田区内幸町1丁目1番6号 日

本電信電話株式会社内

(74)代理人 弁理士 草野 卓

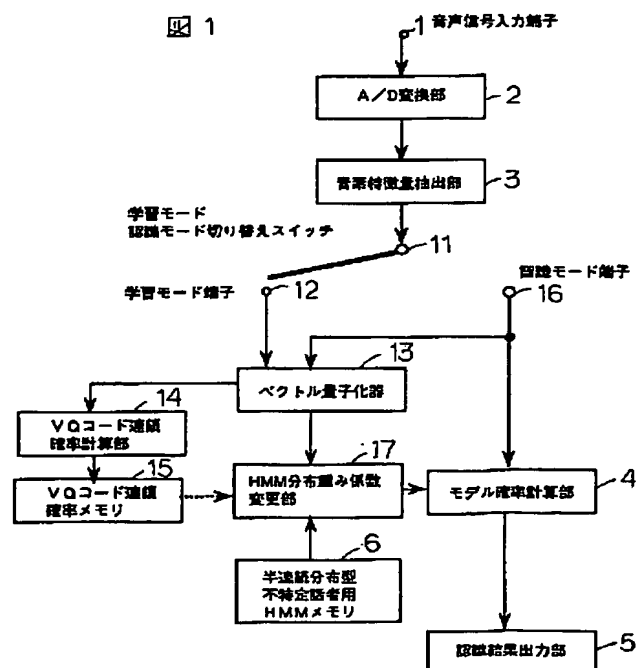
(54)【発明の名称】 不特定話者用音声認識装置

(57)【要約】 (修正有)

【目的】 認識率を向上させる。

【構成】 多数話者の音声进行学习データとして、音素の特徴量を求め、半連続分布型隠れマルコフモデル(HMM)を作成しHMMメモリ6に記憶する。認識に先立ち、入力話者は話者適応化用音声を入力し、それをベクトル量子化し、量子化コードについて、前フレームでコード $c_i$ が出現し、現フレームでコード $c_j$ が出現する連鎖確率 $p(c_j | c_i)$ を計算、各 $p(c_j | c_i)$ をVQコード連鎖メモリ15に蓄える。次に認識すべき単語、の入力音声进行ベクトル量子化しコード列に変換、これと対応した $p(c_j | c_i)$ をメモリ15から読出し、これを用いて、メモリ6から読出したHMMモデルの重み $\lambda_j$ を、 $\lambda_j' = p(c_j | c_i) \lambda_j / \sum (p(c_m | c_i) \lambda_m)$ を演算して変更する。重み係数が変更されたHMMモデルを用いて音素または単語に対する確率を計算し、最大確率と対応する音素又は単語を認識結果として出力する。

図1



1

## 【特許請求の範囲】

【請求項1】 多数話者が発声した音素または単語の音素特微量の分布を、半連続確率密度関数によって表現した隠れマルコフモデルを記憶した不特定話者用隠れマルコフモデルメモリを備え、不特定話者が発声した入力音声から音素特微量を抽出し、その入力音声の発声内容を認識する不特定話者用音声認識装置において、上記入力話者がその不特定話者用音声認識装置を使用する前に発声した音声、または上記入力話者とは異なる多数の話者が発声した音声からそのベクトル量子化コード系列の連鎖確率を計算するベクトル量子化コード連鎖確率計算手段と、上記計算されたベクトル量子化コードの連鎖確率を用いて上記不特定話者用隠れマルコフモデルの半連続確率密度関数の重みを変更する隠れマルコフモデル分布重み係数変更手段と、上記入力話者が発声した認識すべき音声の音素特微量に対し、上記分布の重みを変更した不特定話者用隠れマルコフモデルを用いて音素または単語モデルに対する確率を計算するモデル確率計算手段と、を設けたことを特徴とする不特定話者用音声認識装置。

## 【発明の詳細な説明】

## 【0001】

【産業上の利用分野】この発明は隠れマルコフモデルを用い、音素単位、あるいは単語単位の認識に基づく不特定話者用の音声認識装置に関するものである。不特定話者用の音声認識装置は、話者適応音声を使用しない不特定話者音声認識装置と、話者適応音声を使用する話者適応型不特定話者音声認識装置とがある。話者適応音声を使用しない不特定話者音声認識装置は、入力された音声をただちに認識しなければならないシステム、例えば、音声自動ダイヤルシステムや音声自動券売機などに用いられる。話者適応型不特定話者音声認識装置は、入力話者がシステムを使用する前にいくらかの音声を発声し、この音声を用いてシステムを入力話者に適応化する。例えば、音声ワープロなど、使用者が限定されており、システムが話者に適応化する時間が許されるシステムに使用される。

## 【0002】

【従来の技術】従来における音素単位、あるいは単語単位の認識に基づく不特定話者用の音声認識装置においては、予め、多数話者の音声を学習データとして、それらの音素特微量（例えばケプストラムや振幅）を短区間フレーム（例えば10ms）ごとに求めた後、音素、あるいは単語ごとに音素特微量を統計的な分布をモデル化する。そのモデル化には、統計的な手法の1つである隠れマルコフモデル（Hidden Markov Model、以下HMMと略す。；例えば「確率モデルによる音声認識」電子通信学会、中川聖一著）を用いた手法がある。多次元の音素特微量分布を複数個の連続分布の和

2

で近似し、かつ、全てのモデルの分布の平均値、共分散を共通にしたHMMを半連続分布型HMMと呼ぶ。ただし、各分布に対する重み係数をモデルごとに変えることで、それぞれのカテゴリーの音素特微量の分布を表現する。半連続分布は例えば図2Aに示すようにモデルa、b間で分布の平均値、分散が共通するが、分布の生み係数が異なる。認識時は、入力音声を上記音素特微量に変換した後、モデルごとに半連続確率密度分布を用いて、入力音声に対する各モデルの確率を計算し、最も大きな確率を出力したモデル（音素、単語）を認識結果とする。

【0003】図2Bに、従来の半連続分布型HMMを用いた音声認識装置の構成例を示す。入力端子1から入力された音声は、A/D変換部2においてディジタル信号に変換される。そのディジタル信号は音素特微量抽出部3において音素特微量が抽出され、モデル確率計算部4において、上記の音素、あるいは単語の半連続分布型HMMのパラメータを半連続型不特定話者用HMMメモリ6から読みだし、入力音声に対する各モデルの確率を計算する。その計算された各モデルについての確率中の最も大きな確率を出力する音素、あるいは単語を認識結果として認識結果出力部5より出力する。

## 【0004】

【発明が解決しようとする課題】不特定話者音声認識においては、様々な話者に対応するために、たくさん話者の音声データを用いてモデルを学習する。しかし、話者のバリエーションが増加するに従い、ある話者のある音素の音素特微量分布が、他の話者の異なる音素の音素特微量分布と重なることがしばしば起こる。例えば、話者Aの音素／イ／が、話者Bの音素／エ／に音響的に類似しており、分布がお互いに重なり合っている場合である。これが、不特定話者音声認識における認識誤りの原因の1つとなっていた。

## 【0005】

【課題を解決するための手段】本発明によれば複数フレーム間のベクトル量子化コードの連鎖確率を用いて、不特定話者用半連続型HMMの各分布に対する重み係数を変え、半連続型HMMの分布を入力話者、あるいは認識すべき音素に適するように変更する。手順は学習モードと認識モードに分けられる。学習モードでは、まず、多数話者の音声データを用いて、不特定話者用半連続型HMMを学習する。半連続型HMMの全ての分布の平均値をセントロイド（部分パタン空間の重心）とするベクトル量子化符号帳を作成する。入力話者が予め発声した適応音声、あるいは予め録音された多数話者の音声を上記符号帳を用いてベクトル量子化し、ベクトル量子化コードの連鎖確率を計算する。

【0006】認識モードでは、はじめに入力音声を上記符号帳を用いてベクトル量子化する。前時刻フレームと現時刻フレームのベクトル量子化コードの連鎖確率を用

3

いて、不特定話者用半連続型HMMの各分布に対する重

$$\lambda_j' = p(c_j | c_i) \lambda_j / \sum (p(c_m | c_i) \lambda_m) \dots (1)$$

ここで、 $p(c_j | c_i)$  は前時刻フレームでベクトル量子化コード  $c_i$  が出現し、現時刻フレームでベクトル量子化コード  $c_j$  が出現する連鎖確率であり、 $\lambda_j$  は半連続型HMMの  $j$  番目の分布に対する変更前の重み係数であり、 $\lambda_j'$  は変更後の重み係数であり、 $\Sigma$  は  $m=1$  \*

$$b(y) = \sum \lambda_m b_m(y) \dots (2)$$

$b_m(y)$  は音素特徴量  $y$  の  $m$  番目の分布に対する確率密度を示し、 $\Sigma$  は  $m=1$  から  $M$  までである。

【0008】

【作用】不特定話者用HMMはどのような話者にも対応できるように、多数の話者の音声データから各モデルの確率密度分布を推定する。この分布は一般に広がった分布でモデル間の分布の重なりが大きく誤認識の原因になっていた。この発明では、例えば2フレーム間のベクトル量子化コードの連鎖確率、すなわち音素特徴量の時間遷移情報を用いて、分布の生みを入力話者、あるいは認識すべき音素に適するように変更する。変更された分布は、不特定話者用の分布よりも制約されるので、モデル間の分布の重なりが減少し、音声認識性能が向上する。

【0009】

【実施例】図1にこの発明の実施例を示し、図2Bと対応する部分に同一符号を付けてある。予め多数話者の音声を用いて、音素または単語の半連続分布型HMMを作成し、半連続分布型不特定話者用HMMメモリ6に記憶しておく。はじめに、学習モード認識モード切り替えスイッチ11を学習モード端子12側にする。話者適応型不特定話者音声認識装置の場合は、入力話者の話者適応化用音声を入力して、ベクトル量子化器13でベクトル量子化し、そのベクトル量子化コードについてベクトル量子化コード連鎖確率計算部14で時系列の連鎖確率  $p(c_j | c_i)$  を計算し、その計算結果をベクトル量子化コード連鎖確率メモリ15に蓄える。直ちに認識しなければならない不特定話者音声認識装置の場合は、予め多数の学習話者の音声を入力し、これらについて同様に連鎖確率を計算してベクトル量子化コード連鎖確率メモリ15に蓄えておく。

【0010】次に、学習モード認識モード切り替えスイッチ11を認識モード端子16側にして、入力話者に認識すべき単語を発声してもらう。その入力音声はベクトル量子化器13によってコード列に変換され、ベクトル量子化コード連鎖確率メモリ15よりそのコード列と対

4

み係数を次式により変更する。

$$\lambda_j' = p(c_j | c_i) \lambda_j / \sum (p(c_m | c_i) \lambda_m) \dots (1)$$

\*から  $M$  までであり、 $M$  は分布の総数 (= 符号帳サイズ) である。

【0007】現フレームの音素特徴量  $y$  に対する確率  $b(y)$  を、変更後の重み係数を用いて  $M$  個の連続分布の和によって求める。即ち次式を演算する。

$$b(y) = \sum \lambda_m b_m(y) \dots (2)$$

応した連鎖確率  $p(c_j | c_i)$  を引き出し、その  $p$

10  $(c_j | c_i)$  を用いて、メモリ6から読出した半連続分布型不特定話者用HMMの分布重み係数  $\lambda_j$  をHMM分布重み係数変更部17にて(1)式を演算して  $\lambda_j'$  に変更する。この係数  $\lambda_{jm}$  を用いて(2)式を計算し、これと他のHMMパラメータをもとに入力音声の各HMMに対する確率をモデル確率計算部4にて計算する。

【0011】このようにして計算した確率中の、最大の確率を与えるモデル(音素、単語)を認識結果として認識結果出力部5から出力する。

【0012】

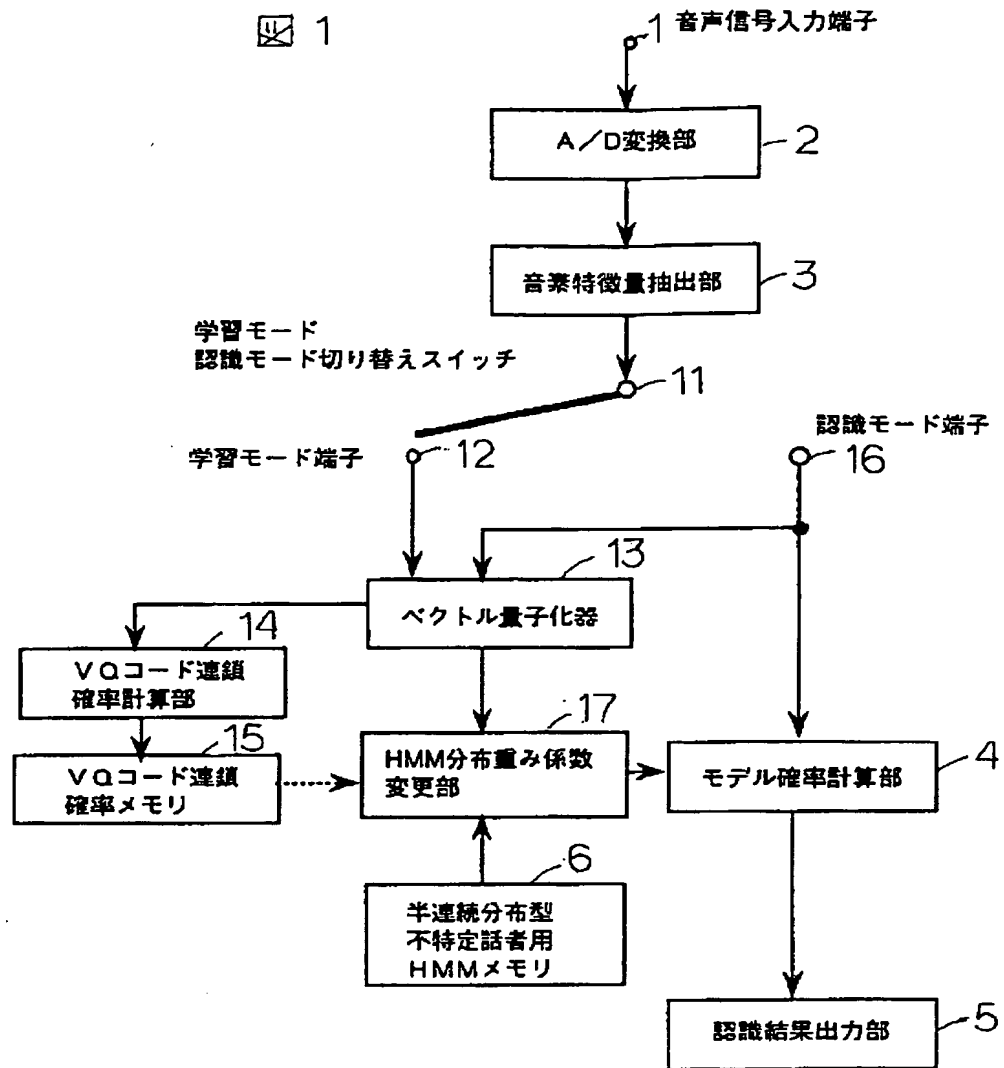
20 【発明の効果】図1に示した構成に従い、不特定話者音声認識を行なった結果を示す。認識対象は連続音声中の23音素である。音素HMMの学習に使用した音声データは64名が発声した9600文章である。評価は上記64名とは異なる13名が発声した823文章である。使用した音素特徴量は16次のケプストラム、16次のデルタケプストラム、1次のデルタパワーである。話者適応型不特定話者音声認識装置として動作させる場合、入力話者のベクトル量子化コード連鎖確率は、各評価話者が発声した50文章より計算した。2つのフレームの時間間隔は8msである。従来の手法では、平均認識率が64.5%であったのに対し、この発明により74.8%にまで改善された。また、直ちに認識する不特定話者音声認識装置として動作させる場合、学習話者64名が発声した9600文章から2000文章を抜き出して、ベクトル量子化コード連鎖確率を計算した。この場合、平均認識率は64.5%から66.5%に改善された。

【図面の簡単な説明】

【図1】この発明の実施例を示すブロック図。

40 【図2】Aは半連続分布型HMMの特徴量の分布例を示す図、Bは従来の不特定話者音声認識装置を示すブロック図である。

【図1】



【図2】

図2 A

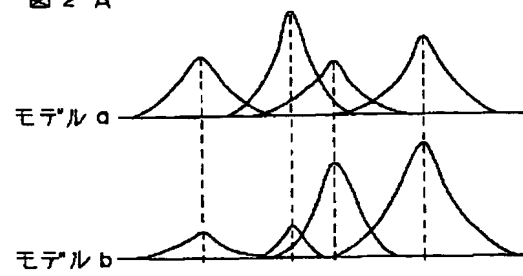
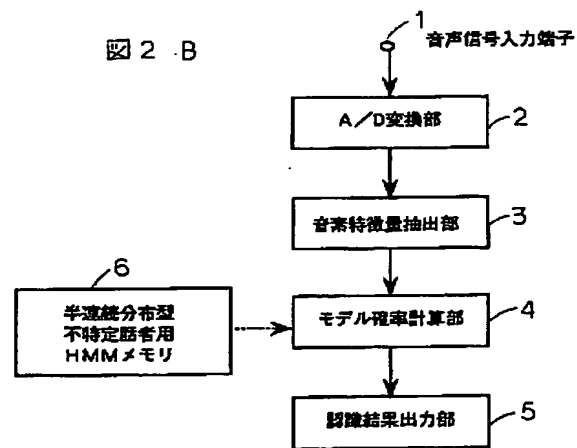


図2 B





**Family list**

**1** family member for:

**JP8241092**

Derived from 1 application.

**1 SPEAKER ADAPTATION METHOD FOR ACOUSTIC MODEL  
AND DEVICE THEREFOR**

Publication Info: **JP8241092 A** - 1996-09-17

---

Data supplied from the **esp@cenet** database - Worldwide

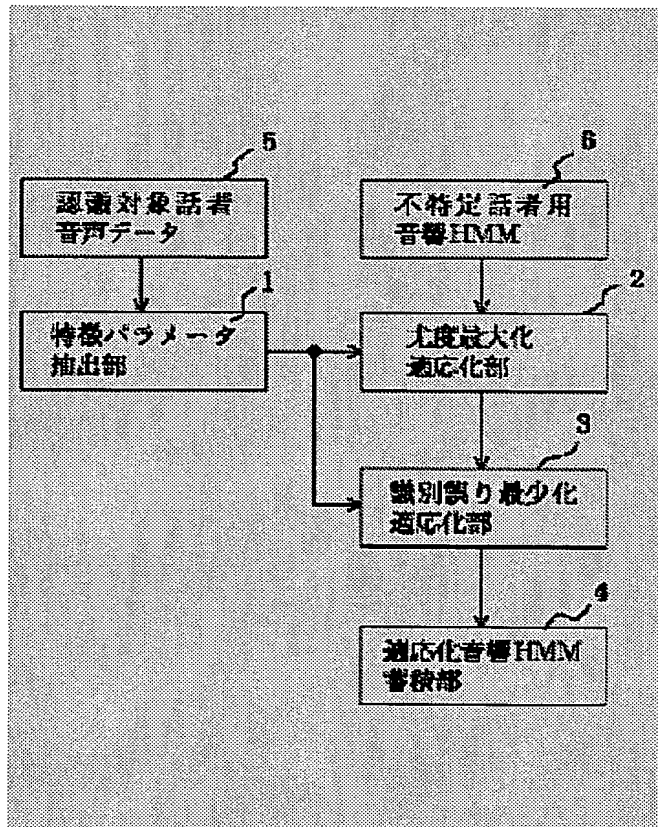
# SPEAKER ADAPTATION METHOD FOR ACOUSTIC MODEL AND DEVICE THEREFOR

**Patent number:** JP8241092  
**Publication date:** 1996-09-17  
**Inventor:** MATSUI TOMOKO; FURUI SADAHIRO  
**Applicant:** NIPPON TELEGR & TELEPH CORP  
<NTT>  
**Classification:**  
- **international:** G10L3/00; G10L3/00; G10L3/00  
- **european:**  
**Application number:** JP19950044430 19950303  
**Priority number(s):**

## Abstract of JP8241092

**PURPOSE:** To realize a recognition system of high performance which minimizes the error rate of an acoustic HMM.

**CONSTITUTION:** This device has a characteristic parameter extraction section 1 which forms and holds the acoustic HMM(hidden Markov model) 6 for nonspecific speakers formed by learning from the speeches of many speakers and extracts a characteristic parameter from speech data 5 of the person to be recognized and a likelihood maximizing adaptation section 2 which optimizes the parameter of the acoustic HMM for the nonspecific persons so as to maximize the likelihood for the speech of a person to be recognized. Further, the device has an identification error minimizing adaptation section 3 which obtains the acoustic HMM of the min. in the identification error by defining the differentiable loss function from the parameter of the acoustic HMM having the parameter maximized in the likelihood and the time series speech data of the characteristic parameter of the person to be recognized and selecting the parameter so as to minimize the function and an adaptation acoustic HMM accumulation section 4.



(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平8-241092

(43) 公開日 平成8年(1996)9月17日

(51) Int.Cl. <sup>6</sup>	識別記号	庁内整理番号	F I	技術表示箇所
G 1 0 L 3/00	5 2 1		G 1 0 L 3/00	5 2 1 F
	5 3 1			5 3 1 K
	5 3 5			5 3 5

審査請求 未請求 請求項の数 4 O L (全 6 頁)

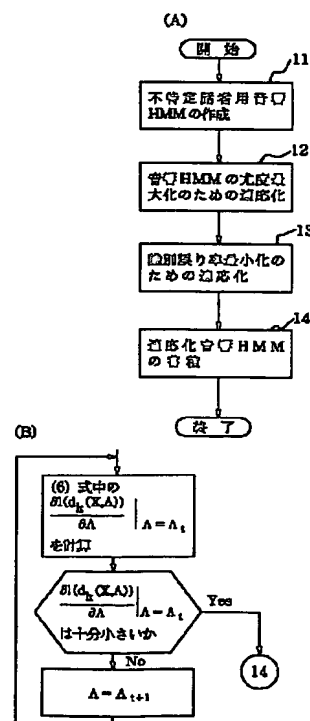
(21) 出願番号	特願平7-44430	(71) 出願人	000004226 日本電信電話株式会社 東京都新宿区西新宿三丁目19番2号
(22) 出願日	平成7年(1995)3月3日	(72) 発明者	松井 知子 東京都千代田区内幸町一丁目1番6号 日 本電信電話株式会社内
		(72) 発明者	古井 貞照 東京都千代田区内幸町一丁目1番6号 日 本電信電話株式会社内
		(74) 代理人	弁理士 若林 忠

(54) 【発明の名称】 音響モデルの話者適応化方法とその装置

(57) 【要約】

【目的】 音響HMMの誤り率を最小化する高性能の音声認識システムの実現。

【構成】 多数の話者の音声から学習して作成した不特定話者用音響HMM (隠れマルコフモデル) 6 を作成して保有し、認識対象話者の音声データ 5 から特徴パラメータを抽出する特徴パラメータ抽出部 1 と、不特定話者用の音響HMMのパラメータを認識対象話者の音声に対する尤度が最大になるように最適化する尤度最大化適応化部 2 と、さらに、尤度を最大にされたパラメータを有する音響HMMのパラメータを認識対象話者の特徴パラメータの時系列音声データとから微分可能な損失関数を定義してその関数が最小になるようにパラメータを選択して識別誤り最小の音響HMMを求める識別誤り最小化適応化部 3 と、適応化音響HMM蓄積部 4 を有する。



1

## 【特許請求の範囲】

【請求項 1】 音声の音響的特徴を抽出し、その特徴量を統計的にモデル化して、音素、単語その他の認識カテゴリに対応した音響モデルを構成するために、多数の話者の音声を用いて学習した不特定話者用の音響モデルを HMM と略称されている隠れマルコフモデルで表現しておき、認識対象となる話者の音声を用いて、前記不特定話者用音響 HMM のパラメータを、認識対象話者の音声に対する尤度が最大となるように最適化する音響モデルの話者適応化方法において、前記認識対象話者の音声に対する尤度が最大になるように最適化された不特定話者用音響 HMM のパラメータを認識対象話者の音声に対する認識誤りが最小になるように適応化するステップを有することを特徴とする音響モデルの話者適応化方法。

【請求項 2】 前記不特定話者用音響 HMM のパラメータを認識対象話者の音声に対する認識誤りを最小になるように適応化するステップが、微分可能な損失関数を定義し、この値が減少するように前記音響 HMM のパラメータを逐次更新して最適値を求めるステップである請求項 1 記載の音響モデルの話者適応方法。

【請求項 3】 音声の音響的特徴を抽出し、その特徴量を統計的にモデル化して、音素、単語その他の認識カテゴリに対応した音響モデルを構成するために、多数の話者の音声を用いて学習した不特定話者用の音響モデルを HMM と略称されている隠れマルコフモデルで表現しておき、認識対象となる話者の音声を用いて、前記不特定話者用音響 HMM のパラメータを、認識対象話者の音声に対する尤度が最大になるように最適化する音響モデルの話者適応化装置において、前記認識対象話者の音声に対する尤度が最大になるように最適化された不特定話者用音響 HMM のパラメータを認識対象話者の音声に対する認識誤りが最小になるように適応化する適応化手段を有することを特徴とする音響モデルの話者適応化装置。

【請求項 4】 前記不特定話者用音響 HMM のパラメータを認識対象者の音声に対する認識誤りを最小になるように適応化する適応化手段が、微分可能な損失関数を定義し、この値が減少するように前記音響 HMM のパラメータを逐次更新して最適値を求める手段を含む請求項 1 記載の音響モデルの話者適応装置。

## 【発明の詳細な説明】

## 【0001】

【産業上の利用分野】本発明は、音声認識方法および装置に関し、特に、音声の音響的特徴量を HMM によってモデル化し、音素、単語などの認識カテゴリに対応した不特定話者用音響 HMM を特定の認識対象話者に適応化する音響 HMM の話者適応化方法と装置に関する。

## 【0002】

【従来の技術】従来、不特定話者音声認識において、認

2

識システムを認識対象話者に適応化することで、認識性能を改善しようとする試みがなされてきた。認識対象話者の音声に対する不特定話者用音響 HMM の尤度が最大になるように、例えば文献「中川聖一：『確率モデルによる音声認識』電子情報通信学会、1988」（以下文献 1 と称す）に発表されたバウム／ウエルチ（Baum-Welch）アルゴリズムに従って、音響 HMM のパラメータを推定して、音 HMM を認識対象話者に適応化していた。

## 【0003】

10 【本発明が解決しようとする課題】上述した従来の音響モデルの話者適応方法は、認識誤り率に関して、良い性能が得られなかったのが実情であった。

【0004】本発明の目的は、認識誤り率に関して、より優れた性能が得られる音響モデルの話者適応化方法とその装置を提供することである。

## 【0005】

20 【課題を解決するための手段】本発明の音響モデルの話者適応化方法は、音声の音響的特徴を抽出し、その特徴量を統計的にモデル化して、音素、単語その他の認識カテゴリに対応した音響モデルを構成するために、多数の話者の音声を用いて学習した不特定話者用の音響モデルを HMM と略称されている隠れマルコフモデルで表現しておき、認識対象となる話者の音声を用いて、前記不特定話者用音響 HMM のパラメータを、認識対象話者の音声に対する尤度が最大となるように最適化する音響モデルの話者適応化方法において、前記認識対象話者の音声に対する尤度が最大になるように最適化された不特定話者用音響 HMM のパラメータを認識対象話者の音声に対する認識誤りが最小になるように適応化するステップを有する。

30 【0006】また、前記不特定話者用音響 HMM のパラメータを認識対象者の音声に対する認識誤りを最小になるように適応化するステップが、微分可能な損失関数を定義しこの値が減少するように前記音響 HMM のパラメータを逐次更新して最適値を求めるステップであるものも本発明に含まれる。

40 【0007】本発明の音響モデルの話者適応化装置は、音声の音響的特徴を抽出し、その特徴量を統計的にモデル化して、音素、単語その他の認識カテゴリに対応した音響モデルを構成するために、多数の話者の音声を用いて学習した不特定話者用の音響モデルを HMM と略称されている隠れマルコフモデルで表現しておき、認識対象となる話者の音声を用いて、前記不特定話者用音響 HMM のパラメータを、認識対象話者の音声に対する尤度が最大となるように最適化する音響モデルの話者適応化装置において、前記認識対象話者の音声に対する尤度が最大になるように最適化された不特定話者用音響 HMM のパラメータを認識対象話者の音声に対する認識誤りが最小になるように適応化する適応化手段を有している。

50 【0008】また、前記不特定話者用音響 HMM のパラ

メータを認識対象者の音声に対する認識誤りを最小になるように適応化する適応化手段が、微分可能な損失関数を定義しこの値が減少するように前記音響HMMのパラメータを逐次更新して最適値を求める手段を含むものも本発明の音響モデルの話者適応装置に含まれる。

【0009】

【作用】多数の話者の音声を用いて、不特定話者用の音響HMMのパラメータを、ある特定の話者の音声に対して尤度が最大になるように最適化した後に、さらに、認識対象話者の音声に対する識別誤りが最小となるように適応化するので、誤り率の少ない音響モデル話者適応化が可能となる。

【0010】

【実施例】次に、本発明の実施例について、図面を参照して説明する。

【0011】図1(A)は本発明の音響モデル話者適応化方法の一実施例のフローチャート、図1(B)は図1\*

$l(d_k(x, \Lambda)) = 1 / [1 + \exp \{-\phi(d_k(x, \Lambda) + \phi)\}]$

を減少するように

【0015】

$$\Lambda_{t+1} = \Lambda_t - \varepsilon_t \frac{\partial l(d_k(x, \Lambda))}{\partial \Lambda}$$

※20

$\Lambda = \Lambda_t$

(6)

ここで $\varepsilon_t$ は更新量を調節する係数で、実験的に設定する。

【0016】式(6)を順次 $\Lambda_t$ を更新して最適値を求めるステップである実施態様も本発明にふくまれる。

【0017】図2は本発明の音響モデル話者適応化方法が適用された装置のブロック図である。

【0018】この音響モデル話者適応化装置は、多数の話者の音声音声を収録して音声の特徴量を統計的にモデル化した不特定話者用の音響モデルを隠れマルコフモデル(以下音響HMM(HIDDEN MARKOV MODEL)と称す)

6と、認識対象話者の音声データ5を入力されると入力された音声をケプストラム等の特徴パラメータを用いた表現形式に変換する特徴パラメータ抽出部1と、特徴パラメータ抽出部1の出力とHMM6が入力されると、認識対象話者の時系列に変換された音声データにより認識対象話者に適応化された尤度の高い音響HMMと、前記認識対象話者の時系列に変換された音声データを出力する尤度最大化適応部2と、尤度最大化適応部2の出力を入力とし、識別誤りを最小にするように音響HMMのパラメータを修正する識別誤り最小化適応部3と、識別誤り最小化適応部3の出力を蓄積する適応化音響HMM蓄積部4を有している。

【0019】尤度最大化適応部2では、入力された音響HMM6のHMMパラメータ $\Lambda =$

【0020】

\* (A) のステップ13の識別誤り率最小化適応化に損失関数を使用した実施態様のフローチャートである。

【0012】この音響モデルの話者適応化方法は、多数の話者の音声を用いて学習した不特定話者用の音響モデルを隠れマルコフモデルである音響HMM(Hidden Markov Model)で表現する(ステップ11)。次に、特定の認識対象話者の音声を用いて音響HMMのパラメータを該認識対象話者の音声に対する尤度が最大になるように最適化する(ステップ12)。さらに、前記認識対象話者の音声に対する識別誤りが最小となるようにステップ12の出力の音響HMMのパラメータを適応化する(ステップ13)。ステップ13の出力である音響HMMを蓄積する(ステップ14)。

【0013】また、識別誤り最小化適応化ステップ13が、損失関数

【0014】

【数1】

※【数2】

【外1】

$\{c_{am}, \mu_{am}, U_{am}, a_{\alpha\beta}\}$

が式1に示す尤度関数 $L_k(\cdot)$ を最大にするように、例えばバウム／ウエルチアルゴリズム(文献1参照)によって認識対象話者に対する適応化を実行する。

【0021】ここで、 $\alpha$ は状態、 $m$ は混合分布、

【0022】

【外2】

$c_{am}$

は状態 $\alpha$ の混合分布 $m$ の重み係数、

【0023】

【外3】

$\mu_{am}$

は平均値、

【0024】

【外4】

$U_{am}$

は分散値、

【0025】

【外5】

$a_{\alpha\beta}$

は状態 $\alpha$ から状態 $\beta$ への遷移確率を表す。

【0026】

【数3】

$$L_K(X, \Lambda) = \sum_Q \pi_{q_0} \cdot \prod_{t=1}^T a_{q_{t-1}, q_t} \cdot \prod_{t=1}^T b_{q_t}(x_t) \quad (1)$$

ここで、 $X = \{x_1, x_2, x_3, \dots, x_T\}$  は特徴パラメータの時系列に変換された音声データ、 $t$  は時刻、 $\pi$  は初期状態確率、 $Q = \{q_0, q_1, \dots, q_T\}$  は状態遷移系列を表す。

\*【0027】また、  
【0028】  
【数4】

$$b_a(x_t) = \sum_{m=1}^M c_{am} \cdot N(x_t | \mu_{am}, U_{am}) \quad (2)$$

は出現確率である。

※ $k(\cdot)$  をそれぞれ式(3)、(4)で定義する。

【0029】識別誤り最小化適応化部3では、識別誤り関数 $d_k(\cdot)$ が最小となるように音響HMMが修正される。ここで、識別関数 $g_k(\cdot)$ と識別誤り関数 $d_k$

【0030】  
【数5】

$$g_k(X, \Lambda) = \log [L_k(X, \Lambda)] \quad (3)$$

$$d_k(X, \Lambda) = -g_k(X, \Lambda) + G_k(X, \Lambda) \quad (4)$$

ここで、

$$G_k(X, \Lambda) = \log \left\{ \left[ \sum \exp [\eta g_j(X, \Lambda)] \right] / (K-1) \right\}^{\frac{1}{\mu}}$$

$L_k$ はビタービ(Viterbi)アルゴリズムによる尤度関数である(文献1による)。

【0031】また、識別誤り最小化適応部3では、例えば音響HMMのパラメータ修正において識別誤り関数 $d_k(\cdot)$ の代りに式5に示す微分可能な損失関数 $l$

☆の音響パラメータを逐次更新する。ここに、式(5)、(6)は文献[B. H. Juang and Katagiri: "Discriminative training" J. Acoust. Soc. Jpn., 13, 6, pp. 333-339, 1992.]の識別学習アルゴリズムによる。

【0032】

( $d_k$ )を定義し、この関数が減少するように式(6) ☆

【数6】

$$l(d_k(x, \Lambda)) = 1 / [1 + \exp \{-\phi(d_k(x, \Lambda) + \phi)\}] \quad (5)$$

$$\Lambda_{t+1} = \Lambda_t - \varepsilon_t \frac{\partial l(d_k(x, \Lambda))}{\partial \Lambda} \quad \Lambda = \Lambda_t \quad (6)$$

ここで、 $\varepsilon_t$ は更新量を調節する係数で、実験的に設定する。

【0033】次に、本実施例の動作結果について実験例を参照して説明する。

【0034】まず、不特定話者用音響HMMの作成を行った。本例の音響HMMは、混合分布数256の半連続型HMMであり、音響HMMは音韻環境独立の43種類である。不特定話者用音響HMMの作成には、男性35名に計7,016文章を用いて、バウムウエルチのアルゴリズムによって、HMMパラメータの推定を行った。

【0035】また、話者適応化に、不特定話者用音響HMMの作成に用いた話者とは異なる男性2名に10および50文章のデータを用いた場合の、連続音声内容の音韻認識率により評価した。音韻認識実験では、音声内容の書き下ろしを与えて音声区間に対してビタビアラインメントを取り、それを正解の音韻区間と仮定し、その音韻区間で全ての音韻HMMのうち、最大尤度を示すものを認識結果とした。音韻識別実験は、話者適応化に用いた

ものと同じ文章セットとは異なる100文章を用いた場合について行った。特徴パラメータとして、標本周波数12KHz、フレーム長32ms、フレーム周期8ms、LPC(Linear Prediction Coefficient)分析次数16でケプストラムを抽出した。

【0036】また、尤度最大化による話者適応化では、学習の繰返し回数は5回とし、また、各混合分布の平均値だけを推定した。また、識別誤り最小化による話者適応化では、繰返し回数は10回として、各混合分布の平均値と重み係数を推定した。両方ともに、各繰返しにおいて、前学習データを適用した後に、一斉にHMMパラメータを更新した。

【0037】上記の実験結果を表1に示したが、本方法が従来の音響モデルの話者適応方法に比して有効であることがよくわかる。

【0038】

【表1】

(数値端音韻認識率%)

話者	学習 10 文音			学習 50 文音		
	適応なし	従来法	本方法	適応なし	従来法	本方法
話者 1	50.1	54.8	58.1	50.1	57.6	70.6
話者 2	36.6	44.2	46.2	36.6	48.4	62.1
平均	43.4	49.5	52.2	43.4	53.0	66.4

## 【0039】

【発明の効果】以上説明したとおり本発明は、不特定話者用音響HMMを尤度最大化適応化された後さらに音響HMMのHMMパラメータの識別誤り最小化適応化をするので、音響モデルの話者適応化が増進され、高性能の音響認識システムを実現できる効果がある。

## 【図面の簡単な説明】

【図1】(A)は本発明の音響モデル話者適応方法の一実施例のフローチャート、(B)は(A)のステップ13の識別誤り率最小化適応化に損失関数を使用した実施

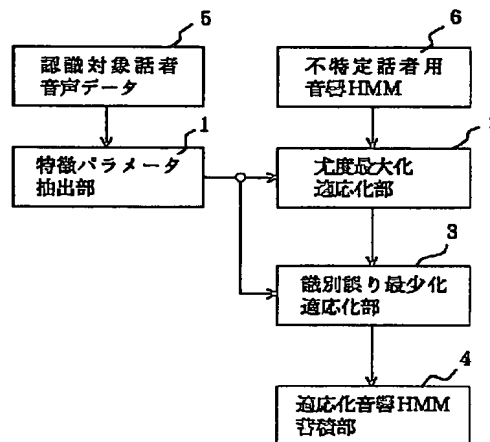
態様のフローチャートである。

【図2】本発明の音響モデルの話者適応装置の一実施例のブロック図である。

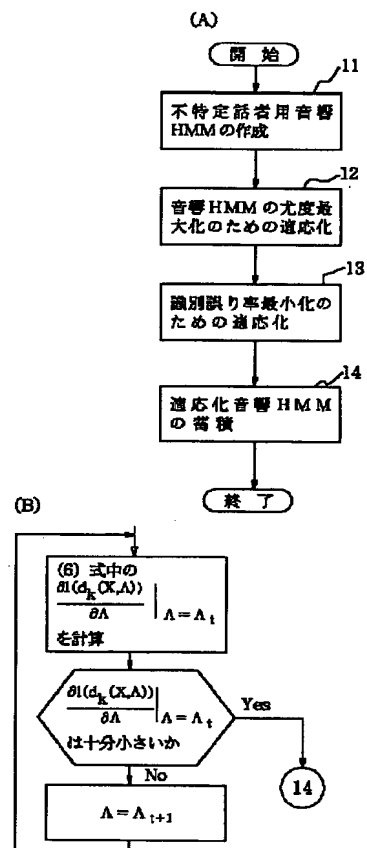
## 【符号の説明】

- 1 特徴パラメータ抽出部
- 2 尤度最大化適応化部
- 3 識別誤り最小化適応化部
- 4 適応化音響HMM蓄積部
- 5 認識対象話者の音声データ
- 6 不特定話者用音響HMM

【図2】



【図 1】





**Family list**

**1** family member for:

**JP62249199**

Derived from 1 application.

**1 SPEAKER CHECKING SYSTEM**

Publication info: **JP62249199 A** - 1987-10-30

---

Data supplied from the *esp@cenet* database - Worldwide

## **SPEAKER CHECKING SYSTEM**

**Patent number:** JP62249199  
**Publication date:** 1987-10-30  
**Inventor:** WATARI MASAO  
**Applicant:** NIPPON ELECTRIC CO  
**Classification:**  
**- international:** G10L3/00  
**- european:**  
**Application number:** JP19860092637 19860421  
**Priority number(s):** JP19860092637 19860421

Abstract not available for JP62249199

---

Data supplied from the **esp@cenet** database - Worldwide

## ⑫ 公開特許公報(A)

昭62-249199

⑤ Int. Cl. 4

識別記号

庁内整理番号

⑬ 公開 昭和62年(1987)10月30日

G 10 L 3/00

3 0 1

D-8221-5D

審査請求 未請求 発明の数 1 (全4頁)

⑭ 発明の名称 話者確認システム

⑮ 特 願 昭61-92637

⑯ 出 願 昭61(1986)4月21日

⑰ 発 明 者 亘 理 誠 夫 東京都港区芝5丁目33番1号 日本電気株式会社内

⑱ 出 願 人 日本電気株式会社 東京都港区芝5丁目33番1号

⑲ 代 理 人 弁理士 内 原 晋

## 明 細 書

発明の名称 話者確認システム

## 特許請求の範囲

あらかじめ用意される複数の抑制標準パターンを記憶する抑制標準パターン記憶部と、登録される確認標準パターンを記憶する手段たる確認標準パターン記憶部と、確認標準パターンに類似した抑制標準パターンを指定する情報を記憶する抑制指定部と、被照合者の発声による入力パターンを保持する入力パターンバッファと、この入力パターンと確認標準パターンとの距離、及びこの入力パターンとこの入力パターンに対応して入力された登録情報指定信号により前記抑制指定部で指定される抑制標準パターンとの距離を計算するパターンマッチング部と、これらの距離が

(イ) 確認標準パターンとの距離 < 第1の閾値  $Th_1$ (ロ) 抑制標準パターンとの距離 > 第2の閾値  $Th_2$ 

の条件を満足したとき被照合者が本人であると決定する判定部とを含んで構成されることを特徴とする話者確認システム。

## 発明の詳細な説明

## (産業上の利用分野)

本発明は話者確認システムの改良に関する。

## (従来技術)

話者確認システムは、音声として入力される合言葉をもつて登録された合言葉の音声パターンと比較して一致の度合を検定することにより、発声者が登録人本人であるかどうかを判定するシステムである。このようなシステムは銀行業務における印鑑にかわるものとして、あるいは入門保安のための錠のかかりとして有用であり、その実現が強く望まれている。

従来試みられている話者確認システムは、例えば日本音響学会誌35巻2号(昭和54年2月)の63頁から69頁に「電話音声を対象とした話者照合」と題して発表された論文に記載される如く、各人のパスワードの発声より作成した標準パターンを記憶して

用いる。このシステムでは、確認動作時には被照合者が主張する名前に対応するコード(以下では簡単に登録番号と呼ぶ)とパスワードの発声が入力される。確認システムでは、入力された登録番号に対応する標準ボタンと、入力されたパスワード音声(以下入力ボタンと呼ぶ)とが比較され距離(ボタン間の相異度の評価値)が算出される。この距離が所定の閾値より小であるときは、この被照合者は本人であると判定され、逆に閾値より大であるときはこの人は詐称者であると判定する。

このような方式の話者確認システムの問題は、閾値の適切な設定が困難であるという点にある。すなわち、同一人が一定のパスワードを発声した場合でも、音声ボタンはその時々で異なっており、しかも変動の程度は個人ごとに異なっている。このため、閾値をきびしく(すなわち小さな値に)設定すると本人が詐称者と判定される事態が多発し、逆に閾値をゆるめに設定すると詐称者を本人であると誤判定してしまうことになる。

(i)抑制標準ボタンとの距離が十分大きいにもかかわらず、本人のボタンのバラツキが大きい場合に(ハ)の条件を満足しないことがあり「他人である」と誤判定されたり、

(ii)抑制標準ボタンとの距離が十分小さいにもかかわらず、他人のボタンのバラツキが大きい場合に(ハ)の条件を満足しているため、「他人を本人である」と誤判定することがあった。

本発明は従来方式の欠点を少なくした、すなわち誤判定の少ない話者確認システムを実現することを目的としている。

(問題を解決するための手段)

本発明による話者確認システムは、あらかじめ用意される複数の抑制標準ボタンを記憶する抑制標準ボタン記憶部と、登録される確認標準ボタンを記憶する手段たる確認標準ボタン記憶部と、確認標準ボタンに類似した抑制標準ボタンを指定する情報を記憶する抑制指定部と、被照合者の発声による入力ボタンを保持する入力ボタンバッファ

この様な欠点の少ない話者確認の方式が特願昭58-054132号明細書に記載されている。その方式の要点は次のとおりである。すなわち、登録される確認標準ボタンの他に、確認標準ボタンに類似する他人のボタンである抑制標準ボタンを用意しておき、入力ボタンとのマッチングの結果得られる距離が次の(イ)、(ハ)条件を満足するときに限って、本人であると判定して一致信号を出す。

(イ)確認標準ボタンとの距離<所定閾値

(ハ)確認標準ボタンとの距離<抑制標準ボタンとの距離

この方式によると(ハ)の判定基準があるために、(イ)で用いる閾値をゆるめに設定できる。これによって本人を詐称者とする誤判定は少なくなり、しかも、他人を本人であるとする誤判定も(ハ)によって防止できるという効果が得られる。

(発明が解決しようとする問題点)

従来方式では、抑制標準ボタンとの距離の大きさの条件がないため、

と、この入力ボタンと確認標準ボタンとの距離、及びこの入力ボタンとこの入力ボタンに対応して入力された登録情報指定信号により前記抑制指定部で指定される抑制標準ボタンとの距離を計算するボタンマッチング部と、これらの距離が

(イ)確認標準ボタンとの距離<第1の閾値Th1

(ロ)抑制標準ボタンとの距離>第2の閾値Th2

の条件を満足したとき被照合者が本人であると決定する判定部とを含んで構成される。

(作用)

本発明の作用について図面を参照しながら説明する。本発明では確認標準ボタンの他に抑制標準ボタンを用意し、入力ボタンとの距離を基に入力ボタンを発声した話者が名のついた本人であるか否かを判定する。ここで確認標準ボタンをA、抑制標準ボタンをB、入力ボタンをX、ボタン間距離をDとすると、判定は次式に従う。

$$D(X, A) < Th1 \quad \text{---(1)}$$

かつ

$$D(X, B) > Th2 \quad \text{---(2)}$$

のとき「本人である」と判定し、それ以外は「他人である」と判定する。すなわち、入力ボタンが第2図に示す斜線部分にある場合「本人である」としていることになる。

(実施例)

第1図は本発明の一実施例を示すブロック図である。抑制標準ボタン記憶部90にはあらかじめ多人数のパスワード音声ボタンが記憶されている。例えば/fujisan/という言葉のボタンが1000人分用意されている。これらの抑制標準ボタンの集合を

$$\{B^1, B^2, \dots, B^m \dots B^{1000}\} \quad \dots(3)$$

とする。各抑制標準ボタン $B^m$ は、例えば特願昭45-84685号明細書(特公昭50-19020号公報)の(1)式に示されるが如く表現されている。

最初に確認標準ボタンの登録に係る動作を説明する。登録番号指定部20はキーボードより成っており、これによってこれから登録する人の情報、すなわち、登録番号 $n$ を入力し、続いて/fujisan/というパスワードを入力する。対応する音声はマイクロホン10より入力される。この音声は分析部

かくして確認標準ボタン $A^n$ が得られ、抑制指定の準備も(4)式の如く整った。以上の処理は新たな被登録人が出現する度に登録番号 $n$ を変えながら繰り返される。

次に確認動作に係る部分を説明する。被照合者は登録番号指定部20のキーを操作して登録情報である登録番号 $n$ を入力するとともに、マイクロホン10より/fujisan/なるパスワード音声を入力する。パスワード音声は、登録時と同様に分析部30によって分析され、入力ボタン $X$ として入力ボタンバッファ40に入力される。

登録番号 $n$ が確認標準ボタン記憶部50に入力されるとそれに応じて確認標準ボタン $A^n$ が出力される。ボタンマッチング部60ではこれを受けて、まず $D(X, A^n)$ なる距離が算出される。

抑制指定部80に登録番号 $n$ が与えられると(4)式の抑制指定信号が信号線 $k$ に出力される。これを受けて抑制標準ボタン記憶部90からは抑制標準ボタン

$$B^m; m=m(n, 1), m(n, 2) \dots m(n, 5) \quad \dots(5)$$

30で分析された後、入力ボタンバッファ40に入力される(以下入力ボタン $X$ と呼ぶ)。この入力ボタン $X$ は確認標準ボタン $A^n$ として登録番号 $n$ に対応づけて確認標準ボタン記憶部50に送られ記憶される。

その後カウンタ100は抑制標準ボタン指定信号 $m$ を1から1000まで順次変化させる。これに応じて抑制標準ボタン記憶部90からは抑制標準ボタン $B^m$ が $B^1$ から $B^{1000}$ まで順次出力され、信号線 $b$ を経由してボタンマッチング部60に送られる。

このボタンマッチング部60では前記の入力ボタン $X$ と、抑制標準ボタン記憶部90から送られてくる抑制標準ボタン $B^m$ との比較処理を行なって距離 $D(X, B^m)$ を順次算出する。抑制指定決定部70では、順次入力される距離 $D(X, B^m)$ が相互に比較され、最小なものから例えば第5番目に小さなものまでが決定され、それに対応する番号 $m$ が抑制指定部80に送られる。抑制指定部80ではこれらの番号 $m$ が、前記登録番号 $n$ に対応づけられて

$$m(n, 1), m(n, 2) \dots m(n, 5) \quad \dots(4)$$

なる形式で記憶される。

が出力される。

ボタンマッチング部60では、前記入力ボタン $X$ と、これら標準ボタン $B^m$ との比較が行なわれ、距離

$$D(X, B^m); m=m(n, 1), m(n, 2) \dots m(n, 5) \quad \dots(6)$$

が算出される。

判定部110では前記距離 $D(X, A)$ と第1の閾値 $Th1$ との比較と(6)式の距離群と第2の閾値 $Th2$ との比較が行われる。この比較の結果

$$D(X, A^n) < Th1 \quad \text{かつ} \quad \dots(7)$$

$$D(X, B^m) > Th2; m=m(n, 1), (n, 2), \dots m(n, 5) \quad \dots(8)$$

のとき、「被照合者が本人である」と決定する。(8)式の判定は、(4)式のごとく $D(X, B^m)$ の最小値を求め、第2の閾値との比較を行っても等価である。すなわち、(8)式は次の(9)式に書きかえることができる。

$$\min_{m=m(n, 1)} D(X, B^m) > Th2 \quad \dots(9)$$

以上本発明の原理を実施例に基づいて説明したが、これらの記載は本発明の権利範囲を限定する

ものではない。特に本明細書では音声パターンを比較するのに距離を用いたが相関のように大小関係が逆の量を用いてもよい。この場合、判定基準(イ)と(ロ)の不等号が逆になるのは自明の理である。

また確認標準パターンを複数個もつことも可能であり、その場合は確認標準パタンのどれかに近ければよい。すなわち(7)式のかわりに(10)式を用いられ

$$\min_j D(X, A_j) < Th1 \quad \dots(10)$$

(発明の効果)

本発明では被照合者が本人であるか否かの判定に前記(イ)、(ロ)の2つの条件を用いているため、本人らしいパターンでも他人である抑制パターンに近い場合「他人である」と判定できる。すなわち、誤判定の少ない話者確認システムを提供することができる。

図面の簡単な説明

第1図は本発明の本実施例を示すブロック図、第2図は本発明の作用を説明する図である。

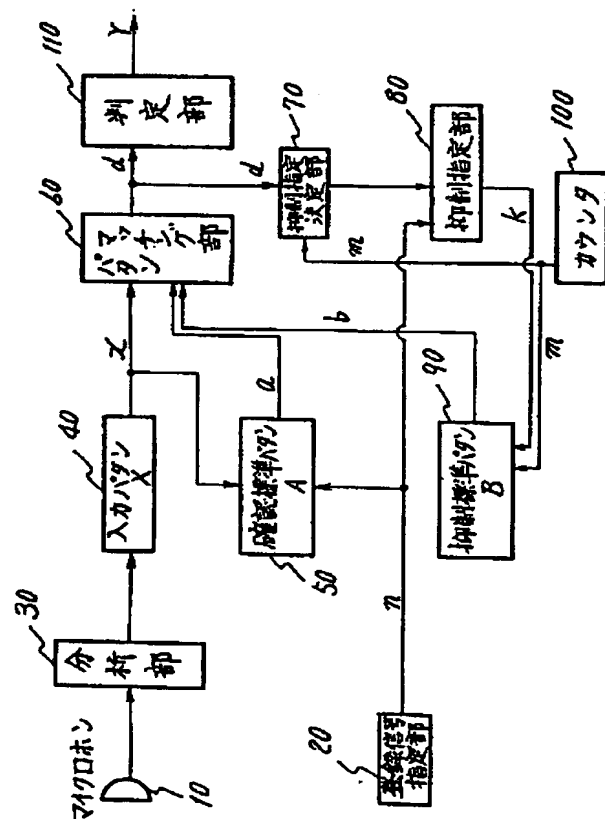
図において、

10…マイクロホン、20…登録番号指定部、30…分析部、40…入力パターンバッファ、50…確認標準パターン記憶部、60…パターンマッチング部、70…抑制指定決定部、80…抑制指定部、90…抑制標準パターン記憶部、100…カウンタ、110…判定部をそれぞれ示す。

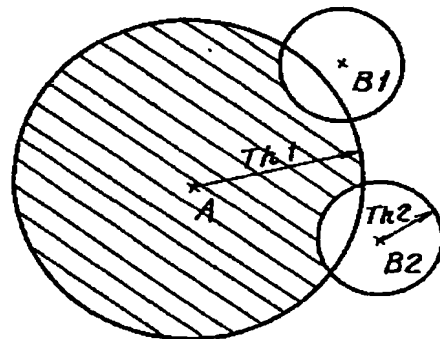
代理人 弁理士 内原 晋



第1図



第2図



**Family list**

**2** family member for:

**JP8123475**

Derived from 1 application.

**1 METHOD AND DEVICE FOR SPEAKER COLLATION**

Publication info: **JP3058569B2 B2** - 2000-07-04

**JP8123475 A** - 1996-05-17

---

Data supplied from the **esp@cenet** database - Worldwide

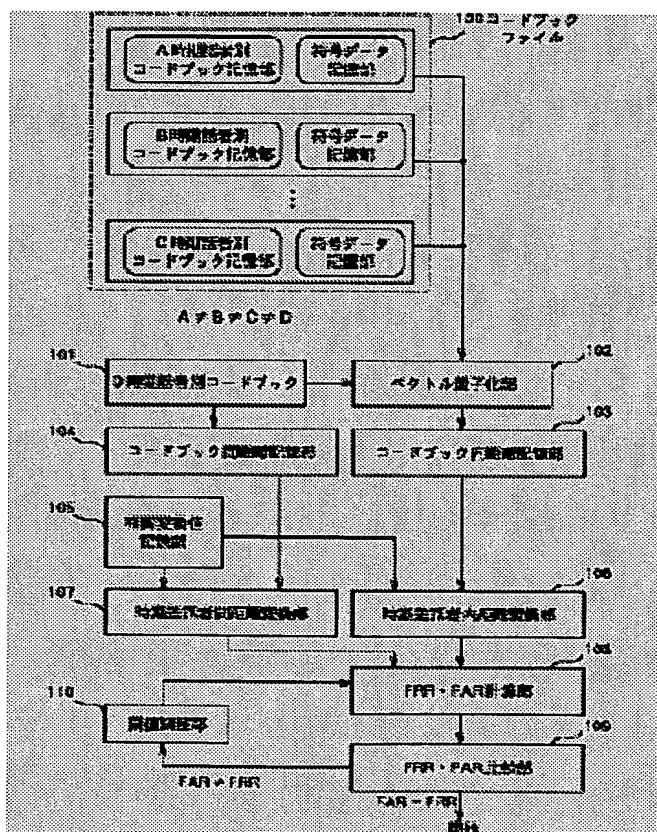
## METHOD AND DEVICE FOR SPEAKER COLLATION

<b>Patent number:</b>	JP8123475
<b>Publication date:</b>	1996-05-17
<b>Inventor:</b>	BIN KATSUTAKE
<b>Applicant:</b>	N T T DATA TSUSHIN KK
<b>Classification:</b>	
- international:	G10L3/00; G10L3/00; G10L9/08
- european:	
<b>Application number:</b>	JP19940265856 19941028
<b>Priority number(s):</b>	

## Abstract of JP8123475

**PURPOSE:** To provide the speaker collation device which considers variance in a speaker and estimates a threshold value adaptive to secular variation of the features of a speech previously in a short time with a small feature sample quantity.

**CONSTITUTION:** Code books by speakers at respective periods are generated on the basis of speech features by the speakers at different periods A-C and stored in a code book file 100. When the threshold value is determined, a code book 101 by the speakers at a period D of a 1st speaker is generated and on the basis of this code book and the past stored code books of the same speaker and other speakers by the speakers, the intra-code-book distance and the inter-code-book distance of the 1st speaker are derived. On the basis of correlative values determined by the speakers at specific time intervals, an intra-period-difference-speaker distance is found from the intra-code-book distance and the inter-period-difference-speaker distance is found from the inter-code-book distance. Then the initial threshold value is adjusted so as to obtain an equal error rate with those intra-speaker distance and inter-speaker distance.





(51) IntCl. <sup>5</sup>	識別記号	庁内整理番号	F I	技術表示箇所
G 1 0 L 3/00	5 3 1 L			
	5 6 1 B			
9/08	3 0 1 C			

審査請求 未請求 請求項の数 5 O L (全 10 頁)

(21) 出願番号 特願平6-265856  
 (22) 出願日 平成6年(1994)10月28日

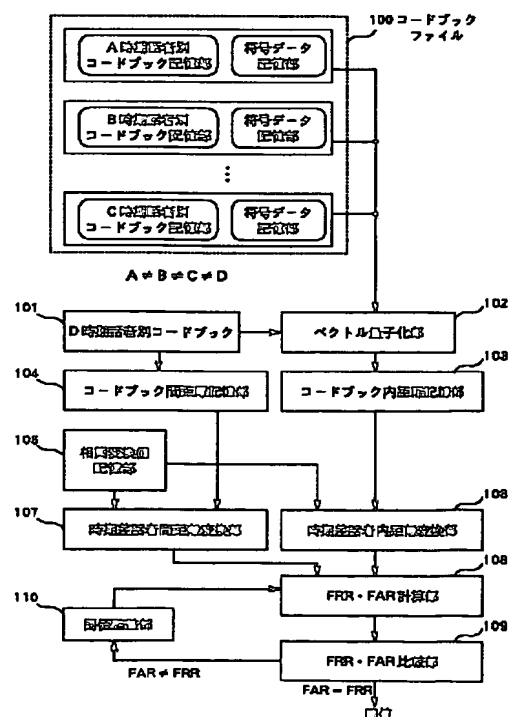
(71) 出願人 000102728  
 エヌ・ティ・ティ・データ通信株式会社  
 東京都江東区豊洲三丁目3番3号  
 (72) 発明者 関 雄偉  
 東京都江東区豊洲三丁目3番3号 エヌ・  
 ティ・ティ・データ通信株式会社内  
 (74) 代理人 弁理士 鈴木 正剛

## (54) 【発明の名称】 話者照合方法及び装置

## (57) 【要約】

【目的】 話者内のばらつきが考慮され且つ声の特徴の経時変化に適応する閾値を少ない特徴サンプル量及び短時間で事前に推定する話者照合装置を提供する。

【構成】 異なる時期A～Cの話者別音声特徴に基づき各時期の話者別コードブックを作成してコードブックファイル100に保存しておく。閾値決定に際しては、まず、第1話者のD時期の話者別コードブック101を作成し、これと保存してある同一話者及び他話者の過去の話者別コードブックに基づき第1話者についてのコードブック内距離、コードブック間距離を導出する。次いで、予め話者別に所定の時期間隔で定めた相関値に基づき、コードブック内距離から時期差話者内距離、コードブック間距離から時期差話者間距離を求める。その後、これら話者内距離、話者間距離による等誤り率になるように初期閾値を調整する。



1

## 【特許請求の範囲】

【請求項 1】 他話者との識別に用いる話者別閾値を決定する閾値決定過程を有する話者照合方法において、前記閾値決定過程は、

前記閾値の決定対象となる第 1 話者の任意の時期の音声特徴に基づき第 1 のコードブックを作成するとともに、この第 1 のコードブックと当該時期以前に作成した第 1 話者のコードブックとの間のコードブック内歪み距離、及び第 1 のコードブックと当該時期以前に作成した他話者のコードブックとの間のコードブック間歪み距離をそれぞれ導出し、更に、第 1 話者について準備された前記コードブック内歪み距離と当該時期の第 1 話者の時期差歪み距離との間の第 1 相関値及びコードブック間歪み距離と他話者の時期差歪み距離との間の第 2 相関値に基づき、当該時期の第 1 話者内の時期差歪み距離及び第 1 話者と他話者との間の時期差話者間歪み距離を導出する過程を含むことを特徴とする話者照合方法。

【請求項 2】 前記閾値決定過程は、更に、当該時期の各時期差歪み距離と任意に定めた初期閾値とに基づき本人拒否率及び詐称者受理率を計算するとともに、これら本人拒否率及び詐称者受理率が等しい値になるように前記初期閾値を調整する過程を含むことを特徴とする請求項 1 記載の話者照合方法。

【請求項 3】 前記第 1 及び第 2 相関値は、予め話者別に所定の時期間隔で取得した値であり、前記第 1 相関値は前記コードブック内歪み距離と同一話者内の時期差歪み距離との間の線形相関値、前記第 2 相関値は前記コードブック間歪み距離と話者間歪み距離との間の線形相関値であることを特徴とする請求項 1 又は 2 記載の話者照合方法。

【請求項 4】 他話者との識別に用いる話者別閾値の決定手段を有する話者照合装置において、各々異なる時期の話者別音声特徴に基づき各時期の話者別コードブックを作成し、各話者別コードブックの作成過程で出現したコードベクトルの出現回数を当該時期の話者別コードブックと共にメモリに保存する手段と、前記閾値の決定対象となる第 1 話者の第 1 のコードブックを作成する対象話者別コードブック作成手段と、作成された第 1 のコードブックと保存してある前記第 1 話者及び他話者の過去のコードブックから各々前記出現回数のコードベクトルを出現させ、これらコードベクトルを前記第 1 のコードブックのコードベクトルで量子化して第 1 話者のコードブック内歪み距離、及び第 1 話者と他話者との間のコードブック間歪み距離を導出する手段と、予め話者別に所定の時期間隔で実測した前記コードブック内歪み距離と同一話者内の時期差歪み距離との間の第 1 相関値、及び前記コードブック間歪み距離と他話者間の時期差歪み距離との間の第 2 相関値を記憶した相関値記憶手段と、

2

前記第 1 話者に対応する前記第 1 及び第 2 相関値を読み出して当該時期の第 1 話者内の時期差歪み距離及び他話者間の時期差歪み距離を導出する時期差歪み距離導出手段と、

を有することを特徴とする話者照合装置。

【請求項 5】 請求項 1 記載の話者照合装置において、各時期差歪み距離と任意に定めた初期閾値とに基づき本人拒否率及び詐称者受理率を計算するとともに、これら本人拒否率及び詐称者受理率が等しい値になるように前記初期閾値を調整する閾値調整手段と、

を有することを特徴とする話者照合装置。

## 【発明の詳細な説明】

## 【0001】

【産業上の利用分野】 本発明は、話者認識システムや音声認識システムに用いられる話者照合技術に関し、特に、発話者の声の特徴の経時変化に適応する話者別閾値を少ない特徴サンプル量で短時間に推定する方法及びその実現装置に関する。

## 【0002】

【従来の技術】 話者照合装置は、発話者の表明した識別名称が、発話者自身の真の識別名称と一致するかどうかを判定する装置である。通常、話者照合を行う場合には、予め照合対象となる話者識別名称及びこの識別名称に対応するコードブックを登録しておき、話者照合時に、発話者の実音声と識別名称とを入力し、この識別名称によって指定されたコードブックと発話者の実音声とを比較してその特徴差を検出する。この特徴差が話者別に設定された所定の閾値以下の場合には表明された識別名称が真の識別名称であり、発話者は本人であると判定する。その他の場合は、表明された識別名称は偽識別名称であり、発話者は詐称者であると判定する。このように、話者照合においては、話者別閾値をどのような値に決定するかが重要であり、この値が適切な値であるかどうかによって話者の識別率が大きく左右される。

【0003】 話者照合時の誤認識には大別して 2 つの原因がある。1 つは発話者が真の識別名称を表明しているにも拘わらず、識別名称が偽であると認識してしまう場合であり、この誤認識率を本人拒否率 (FRR: False Rejection Rate) と称する。もう 1 つは、発話者が偽名称を表明しているにも拘わらずそれを真の識別名称と認識してしまう場合であり、この誤認識率を詐称者受理率 (FAR: False Acceptance Rate) と称する。ところで、話者別閾値の値を高くすると、特徴差が大きくても発話者が本人であると判断する確率が高くなる。従って、FRR は低くなるが FAR は高くなる。逆に、話者別閾値の値を低くすると、FAR は低くなるが FRR は高くなる。このように、FRR と FAR とは一方が低くなると他方が高くなるという関係にある。誤認識率は両者の平均値で表されるので、話者別閾値を調整して両者の平均値をできるだけ小さくすることが好ましい。

【0004】従来、この閾値を決定するための手法が種々提案されている。第1の手法として、FRRとFARとが等しくなるように話者別閾値を設定する等誤り率設定法があり、“デジタル音声処理”（著者：古井 貞照；出版者：東海大学出版会）第9章に紹介されている。図4は、この等誤り率設定法を実現するためのブロック図であり、本人学習音声及び詐称者学習音声を音声入力端子400に入力し、前処理部401が各音声を一定時間長の音声フレームごとに記憶する。特徴量抽出部402は、各音声の特徴量を抽出する。ベクトル量子化部403は、音声から抽出されたそれぞれの特徴量を識別名称に対応する話者コードブック404に基づいてベクトル量子化し、これにより得られたコードベクトルの同一話者内歪み距離（以下、話者内距離）、他話者間歪み距離（以下、話者間距離）を話者内／話者間距離記憶部405に記憶する。FRR・FAR計算部406は、話者内距離と予め定められた初期閾値とを用いてFRRを計算するとともに、話者間距離と上記初期閾値とを用いてFARを計算する。FRR・FAR比較部408では、FRRとFARの値を比較し、両者が等しくなければ閾値調整部407において初期閾値を調整し、再度FRR・FAR計算部406に戻る。そしてFRRとFARとが等しくなった時点で、調整を終え、その値を当該話者の閾値として出力する。

【0005】また、第2の手法として、話者間距離の分布を考慮して閾値を設定する方法（S.Furui, “Cepstral Analysis Technique for Automatic Speaker Verification,” IEEE Trans. Acoustics, Speech, and Signal Processing, vol. ASSP-29, No. 2, pp. 254-272, April 1981 参照）が知られている。図5はこの手法を実現するためのブロック図であり、音声入力端子500に学習音声が入力された後、前処理部501、特徴量抽出部502、ベクトル量子化部503、話者コードブック504までは図4の構成と同様となる。この手法の特徴は、ベクトル量子化部503で得られた話者間距離を話者間距離記憶部505に記憶しておき、それぞれ話者間標準偏差計算部506と話者間平均値計算部507において話者間距離の平均値と標準偏差を求め、その結果得られた統計パラメータに基づき閾値計算部508で閾値を導出することにある。

【0006】また、第3の手法として、本発明者らにより開示された「話者照合方法及び装置」（特願平6-41615号明細書参照）がある。この手法は、複数のコードブックから発話者の表明した識別名称に対応する本人コードブックとそれ以外の他話者コードブックとを選択し、他話者コードブックから出現した所定量のコードベクトルと本人コードブックとの特徴差の統計値を計算し、これにより閾値を得るものである。つまりコードブック間距離を話者間距離に変換することを特徴とする。この手法は、図6の各ブロック600～612により実

現される。

#### 【0007】

【発明が解決しようとする課題】上記各手法は、いずれも特定の一時期に収集した特徴サンプルに基づいて閾値を決定する手法であり、人間の声の特徴に経時変化があることを考慮していない。そのため、時間が経つにつれて話者の識別率が低下する場合があった。人間の声の特徴の経時変化に適応する閾値を推定するには、発話者毎にできるだけ長期間の特徴サンプルを用いて音声特徴の標準パターンを作成しておくことでその対処が可能である。しかしながら、長期間の特徴サンプルをそのまま保存する場合或いは音声特徴を抽出して保存する場合のいずれであっても、話者照合装置に莫大なメモリ容量を確保しなければならず、しかも、特徴サンプル等が膨大な量になることから話者別閾値の計算に長時間を要する問題があった。

【0008】また、上記第1及び第2の手法は、あくまでも事後的に閾値を設定する手法なので、推定等によって事前に閾値の設定を要する用途では十分に活用できず、また、第2及び第3の手法は、話者内距離のばらつきを考慮しないため、話者照合時に、本人を高い確率で拒否してしまう可能性があった。

【0009】本発明の課題は、上記問題点を解消し、話者内距離のばらつきが考慮され、しかも声の特徴の経時変化に適応する閾値を、少ない特徴サンプル量及び短時間で事前に推定する方法及びこの方法を実施するための装置を提供することにある。

#### 【0010】

【課題を解決するための手段】本発明は、コードブックサイズが一定値以上であれば、コードブック内距離と話者内距離、コードブック間距離と話者間距離との間に、それぞれ図2及び図3に示すように強い相関関係があり、しかもこれらの相関関係は時期差に頑健であるという性質を有効に利用して、話者内距離のばらつき及び同一話者及び他話者の声の特徴の時期差を考慮した最適な話者別閾値を決定する点に特徴がある。

【0011】図2は、例えば所定語数から成る単位大きさの文章の音声を所定フレーム長で抽出した特徴量を1セットとした場合に、8セットの男性音声に基づくコードブック内距離と話者内距離との間の相関実測図である。例えばコードブックサイズは512であり、時期差は9ヶ月である。この図から明らかなように、同一話者であれば両者は線形相関にあり、時期差話者内距離を $y$ 、コードブック内距離を $x$ とすると、 $y = ax + b$ の関係にあることが本発明者による検証の結果明らかになった。この式において係数 $a$ 、 $b$ については多少のばらつきはあるものの、全体的には時期差に頑健な傾向が現れている。なお、図2の例では、 $a$ は0.944、 $b$ は0.101であった。

【0012】また、図3は、1セットの女性音声に基づ

くコードブック間距離と時期差話者間距離との間の相関実測図であり、図 2 の場合と同様、コードブックサイズは 512、時期差は 9 ヶ月である。図 3 を参照すると 1 セットであるにも拘わらず、図 2 と同様の線形相関であることが明らかであり、しかもこの傾向は、セット数が増えても同様となる。なお、図示を省略したが、男性音声に基づくコードブック間距離と話者間距離との関係も同様であった。即ち、コードブック内距離とコードブック間距離が判れば、これら相関関係に基づき、時期差話者内距離と時期差話者間距離を導出することができる。

【0013】このような性質を利用した本発明の話者照合方法は、閾値の決定対象となる第 1 話者の任意の時期の音声特徴に基づき第 1 のコードブックを作成するとともに、第 1 話者のコードブック内距離及びコードブック間歪み距離をそれぞれ導出し、更に、第 1 話者について準備された前記コードブック内距離と当該時期の時期差話者内距離との間の第 1 相関値及びコードブック間距離と時期差話者間距離との間の第 2 相関値に基づき、当該時期における時期差話者内距離及び時期差話者間距離を導出する過程を経る。時期差話者内距離及び話者間距離を導出した後は、従来の第 1 手法と同様に、これら各距離と任意に定めた初期閾値とに基づき本人拒否率及び詐称者受理率を計算するとともに、これら本人拒否率及び詐称者受理率が等しい値になるように前記初期閾値を調整すれば良い。

【0014】なお、第 1 及び第 2 相関値は、図 2 及び図 3 から明らかなように、予め話者別に所定の時期間隔で取得した値であり、第 1 相関値は前記コードブック内歪み距離と同一話者内の時期差歪み距離との間の線形相関値、第 2 相関値は前記コードブック間歪み距離と話者間歪み距離との間の線形相関値である。

【0015】また、上記性質を利用した本発明の話者照合装置は、各々異なる時期の話者別音声特徴に基づき各時期の話者別コードブックを作成し、各話者別コードブックの作成過程で出現したコードベクトルの出現回数を当該時期の話者別コードブックと共にメモリに保存する手段を有する。この手段は、公知技術を利用することで実現することができる。また、閾値の決定対象となる第 1 話者の第 1 のコードブックを作成する対象話者別コードブック作成手段と、作成された第 1 のコードブックと保存してある前記第 1 話者及び他話者の過去のコードブックから各々前記出現回数のコードベクトルを出現させ、これらコードベクトルを前記第 1 のコードブックのコードベクトルで量子化して第 1 話者のコードブック内歪み距離、及び第 1 話者と他話者との間のコードブック間歪み距離を導出する手段と、予め話者別に所定の時期間隔で実測した前記コードブック内歪み距離と同一話者内の時期差歪み距離との間の第 1 相関値、及び前記コードブック間歪み距離と他話者間の時期差歪み距離との間の第 2 相関値を記憶した相関値記憶手段と、前記第 1 話

者に対応する前記第 1 及び第 2 相関値を読み出して当該時期の第 1 話者内の時期差歪み距離及び他話者間の時期差歪み距離を導出する時期差歪み距離導出手段と、各時期差歪み距離と任意に定めた初期閾値とに基づき本人拒否率及び詐称者受理率を計算するとともに、これら本人拒否率及び詐称者受理率が等しい値になるように前記初期閾値を調整する閾値調整手段と、を有する。

【0016】

【作用】本発明では、長期間の特徴サンプルの代わりに、話者別に異なる時期の話者別コードブックを複数作成しておき、その際、各話者別コードブックにおけるコードベクトルの出現回数を保存しておく。また、好ましくは閾値決定前に図 2 及び図 3 で示した相関関係、即ち第 1 相関値及び第 2 相関値を相関値記憶手段に記憶させておく。第 1 話者の閾値を決定するときは、過去の各時期の第 1 話者及び他話者のコードブックからコードベクトルを代表する符号列及び符号列の出現回数に従ってコードベクトルを出現させ、閾値決定時期に対応するコードブック内距離、及びコードブック間距離を求める。次いで、このコードブック内距離、コードブック間距離と上記相関値記憶手段から読み出した第 1 相関値及び第 2 相関値に基づき、時期差話者内距離、時期差話者間距離を導出する。

【0017】このように過去において話者別に作成されたコードブックを用いるだけで事前に第 1 話者の時期差話者内距離と時期差話者間距離を近似表現することが可能となり、話者照合装置におけるメモリ使用量が従来よりも大幅に節約される。因みに従来の各手法により音声波形をそのまま保存する場合 (short 型) は、サンプリング周波数 ( $1/\text{sec}$ )  $\times$  音声継続時間 (sec)  $\times 2$  (バイト)、音声波形からその特徴ベクトルを抽出して保存する場合 (double 型) は、総フレーム数  $\times$  特徴ベクトルサイズ  $\times 16$  (バイト) のメモリ容量を必要とするのに対し、本発明の方法及び装置の場合のメモリ使用量は、話者別コードブックサイズ  $\times$  特徴ベクトルサイズ  $\times 16$  (バイト: double 型) + 話者別コードブックサイズ  $\times 4$  (バイト: int 型) である。従って、サンプリング周波数が 16 KHz の単語音声 10 個あるとし、その平均の長さが 3 秒、分析フレーム長が 30 msec、フレーム周期が 16 msec とすると、1 ヶ月毎に特徴サンプルの収録を重ね、1 年間で収録した 100 名の話者の特徴サンプルの特徴量のみを記憶するために必要なメモリ容量は、コードブックサイズが 256 であれば従来の  $1/11$  倍となる。つまり、本発明によれば、約 92% のメモリ容量が節約できる。計算量についても同様であり、本発明によれば、従来の約 92% の計算量が削減できる。上記時期差話者内距離と時期差話者間距離が導出された後は、従来の第 1 の手法と同様の手順乃至手段で等誤り率になるように初期閾値を調整し、これを第 1 話者の閾値とする。

## 【0018】

【実施例】次に、図面を参照して本発明の実施例を詳細に説明する。図1は、本発明の一実施例の話者照合装置における話者別閾値決定部のブロック図であり、100はコードブックファイル、101は話者別コードブック（D時期話者別コードブック）、102はベクトル量子化部、103はコードブック内距離記憶部、104はコードブック間距離記憶部、105は相関変換値記憶部、106は時期差話者内距離変換部、107は時期差話者間距離変換部、108はFRR・FAR計算部、109はFRR・FAR比較部、110は閾値調整部である。

【0019】コードブックファイル100には、話者別に異なる時期、図示の例ではA～C時期に作成したコードブックを記憶する話者別コードブック記憶部と、各話者別コードブックのコードベクトルを代表する符号列及び各符号の出現回数を記憶する符号データ記憶部とが格納され、話者別の過去の音声特徴として随時再利用できるようにになっている。時期A～Cの間隔は、ある程度離れた方が好ましい。これは、短期間では話者の音声特徴に差異が生じない場合があるからである。また、一回も出現しなかった符号乃至符号列については符号データ記憶部への記憶を行わない。これによってコードブックファイル100のサイズ（メモリ使用量）及び以後の計算量の削減が更に期待できる。なお、符号列とは、各コードベクトルに対し、例えばそれぞれのクラスタのセントロイドに対応して付与された符号の集合をいい、符号の出現回数とは、ベクトル量子化処理が終了するまでの過程において、同じクラスタに配属された符号の出現回数データをいう。

【0020】D時期の話者別コードブック101は、閾値設定対象となる第1話者（甲）のD時期における学習音声から分析フレームを抽出して特徴量を求め、その特徴量をベクトル量子化して作成したコードブックである。このD時期は、任意の時期であるが、上述のA時期、B時期、C時期よりも遅い時期である。

【0021】ベクトル量子化部102は、例えば上記D時期話者別コードブック101に基づき、コードブックファイル100内に既に格納されている同一話者（甲）のA～C時期のコードブックからコードベクトルを各符号の出現回数に従って出現させ、これを例えばD時期話者別コードブック101に出力させる。これにより同一話者（甲）によるD時期のコードブックの特徴差、即ちコードブック内距離が得られる。これをコードブック内距離記憶部103に記憶する。

【0022】また、コードブックファイル100内に既に格納されている他話者（乙）のA～C時期の話者別コードブックからコードベクトルを各符号の出現回数に従って出現させ、これを例えばD時期の話者別コードブック101に出力させる。これにより上記話者と他話者によるD時期のコードブック特徴差、即ちコードブック間

距離が得られる。これをコードブック間距離記憶部104に記憶する。

【0023】相関変換値記憶部105には、話者別に、図2及び図3に基づく線形相関式とコードブック内距離／コードブック間距離から一義的に導かれる相関値が記憶されている（相関値記憶手段）。ここにいう相関値は、コードブック内距離と時期差話者内距離との関係を表す第1相関値、コードブック間距離と時期差話者間距離との関係を表す第2相関値である。これら相関値は、各話者をキーに読み出され、各々第1相関値は時期差話者内距離変換部106、第2相関値は時期差話者間距離変換部107に入力される。この実施例の場合は、話者

（甲）及び他話者（乙）に関わる第1及び第2相関値をそれぞれ読み出す。

【0024】時期差話者内変換部106及び時期差話者間距離変換部107は、それぞれ上記コードブック内距離及びコードブック間距離と上記第1及び第2相関値に基づいてD時期における話者（甲）の時期差話者内距離及び時期差話者間距離を求め、これらをFRR・FAR計算部108に入力する。

【0025】FRR・FAR計算部108は、入力された時期差話者内距離と任意の初期閾値とを用いてFRRを計算するとともに、入力された時期差話者間距離と初期閾値とを用いてFARを計算する。FRR・FAR比較部109では、FRRとFARの値を比較し、両者が等しくなければ閾値調整部110において初期閾値を調整し、再度FRR・FAR計算部108に戻る。そしてFRRとFARとが等しくなった時点で、調整を終え、そのときの閾値を当該話者（甲）の閾値として出力する。

【0026】このように、本実施例では、予めコードブック内距離と時期差話者内距離との相関関係、及びコードブック間距離と時期差話者間距離との相関関係を求めて相関変換値記憶部104に記憶しておき、また、話者別に異なる時期に作成した複数の話者別コードブックを保存しておき、閾値決定の際に、これら話者別コードブックを用いてコードブック内距離及びコードブック間距離を導出するとともに上記相関関係に基づき時期差話者内距離及び時期差話者間距離を導出するようにしたので、最初の時期（コードブック作成時）を除けば以後の各時期において話者別閾値を随時決定することができる。また、コードブックファイル100のメモリ使用量が従来手法に比べて大幅に節約され、閾値計算に要する時間も短縮化される。これにより従来の課題を解決することができる。なお、本実施例では、第1及び第2相関値を予め相関変換値記憶部105に記憶しておき、閾値決定対象となる話者をキーとして読み出す構成としたが、必ずしもこのような構成に限定されるものではなく、随時計算して時期差話者内距離変換部106及び時期差話者間距離変換部107に入力するようにしても良

い。

# 【0027】

【発明の効果】以上の説明から明らかなように、本発明の話者照合方法によれば、コードブック内距離及びコードブック間距離と予め取得した第1相関値及び第2相関値に基づいて時期差話者内距離及び時期差話者間距離が導出されるので、声の特徴の経時変化に対応する話者別閾値を決定する際に、長期間の特徴サンプルを記憶する必要がなくなり、話者照合装置のメモリ使用量を節約することができる。更に、特徴サンプル量が少なくて済むことから計算時間が従来手法よりも大幅に短縮され、話者別閾値を、短時間で推定し得る効果がある。

【0028】また、本発明の話者照合装置によれば、過去の異なる時期に話者別に作成したコードブックを再利用して事前に第1話者の時期差話者内距離と時期差話者間距離を相関値によって近似表現することが可能となるので、声の特徴の経時変化に適応する話者別閾値を決定する場合であってもメモリを余分に使用する必要がなくなる。また、特徴サンプル量が少なく、各距離の計算時間も短縮化されるので、話者別閾値の決定手段の全体構成を簡略化し得る効果がある。

## 【図面の簡単な説明】

【図1】本発明の一実施例に係る話者照合装置の要部ブロック図。

【図2】コードブック内距離と話者内距離との相関関係を示す実測図。

【図3】コードブック間距離と話者間距離との相関関係を示す実測図。

【図4】従来の第1の手法である等誤り閾値設定方法を実現するためのブロック図。

【図5】従来の第2の手法である話者間距離の分布を考慮した閾値設定方法を実現するためのブロック図。

【図6】従来の第3の手法であるコードブック間距離を話者間距離に変換する方法を実現するためのブロック図。

## 【符号の説明】

100 コードブックファイル

101 任意の時期、例えばD時期に作成した話者別コードブック

102 ベクトル量子化部

103 コードブック内距離記憶部

104 コードブック間距離記憶部

105 相関変換値記憶部（相関値記憶手段）

106 時期差話者内距離変換部

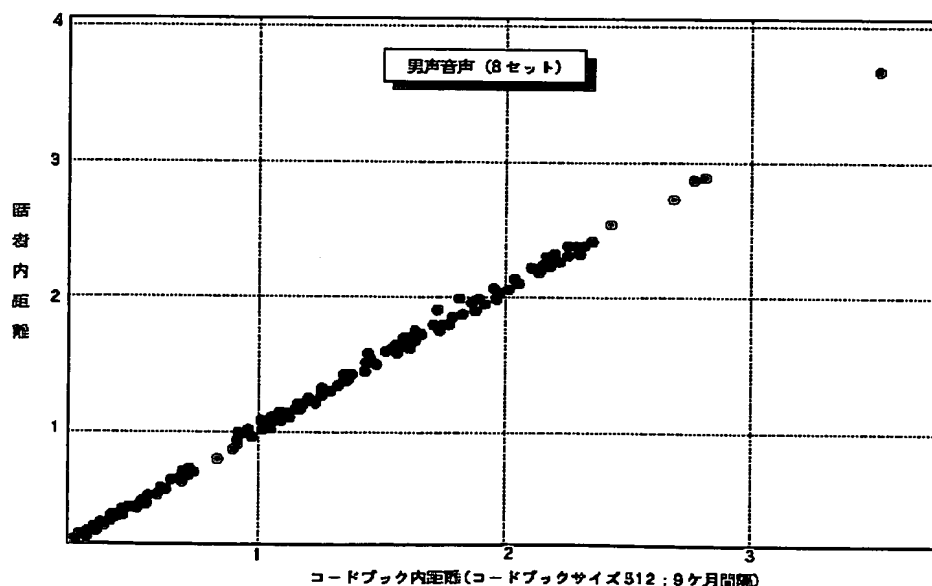
107 時期差話者間距離変換部

108 FRR・FAR計算部

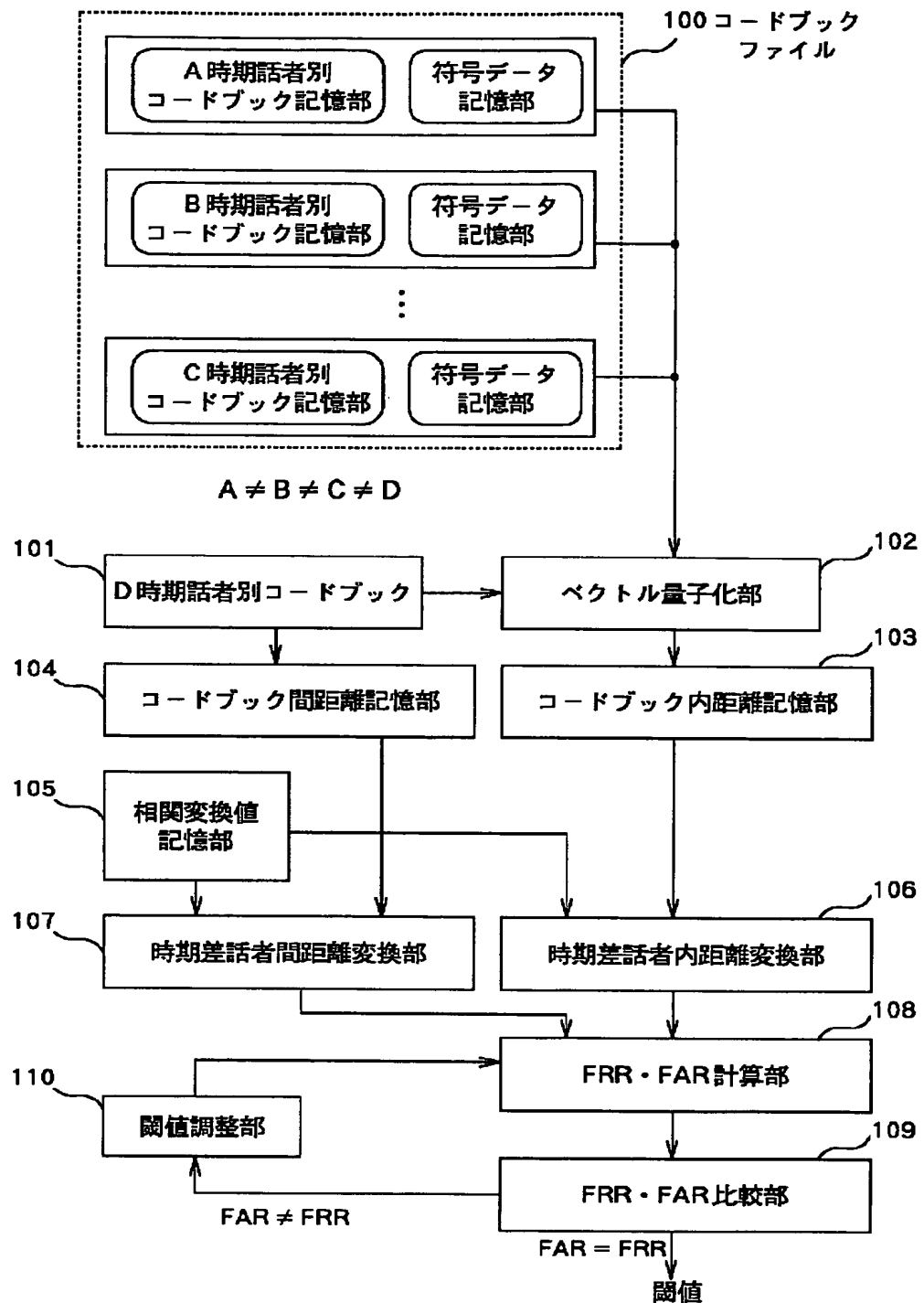
109 FRR・FAR比較部

110 閾値調整部

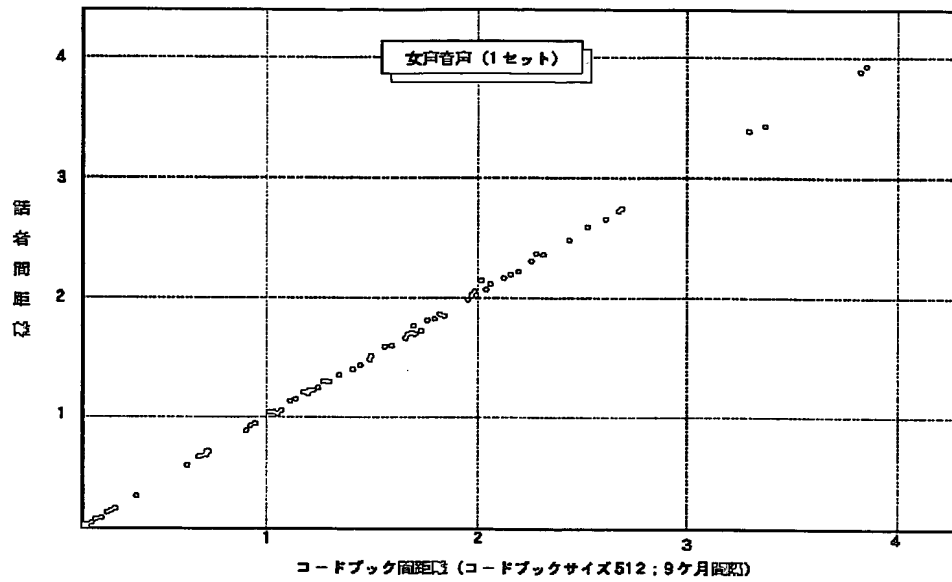
【図2】



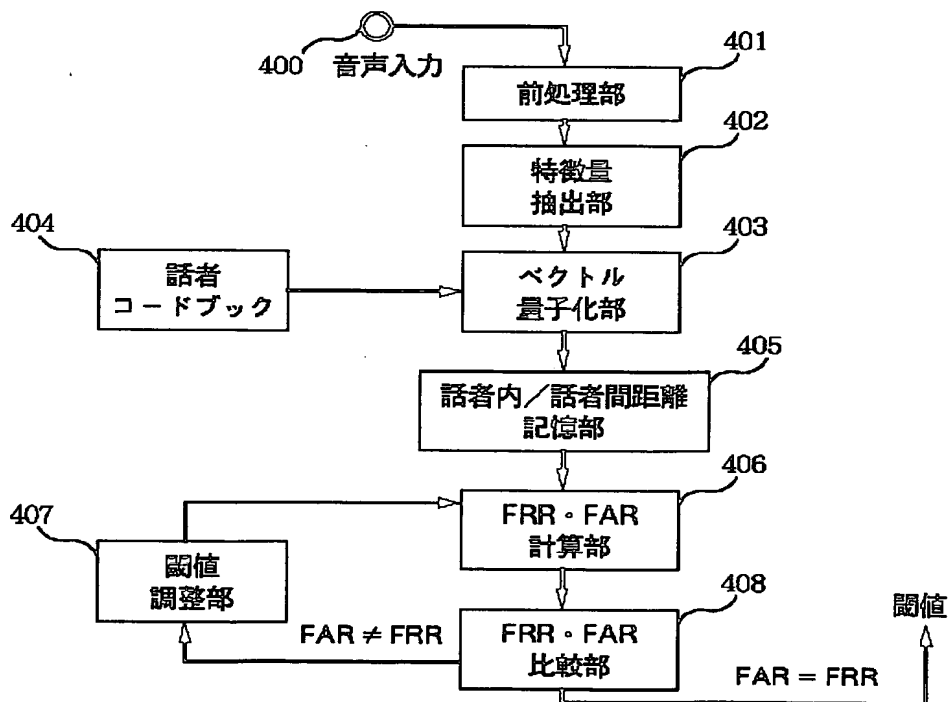
【図1】



【図 3】

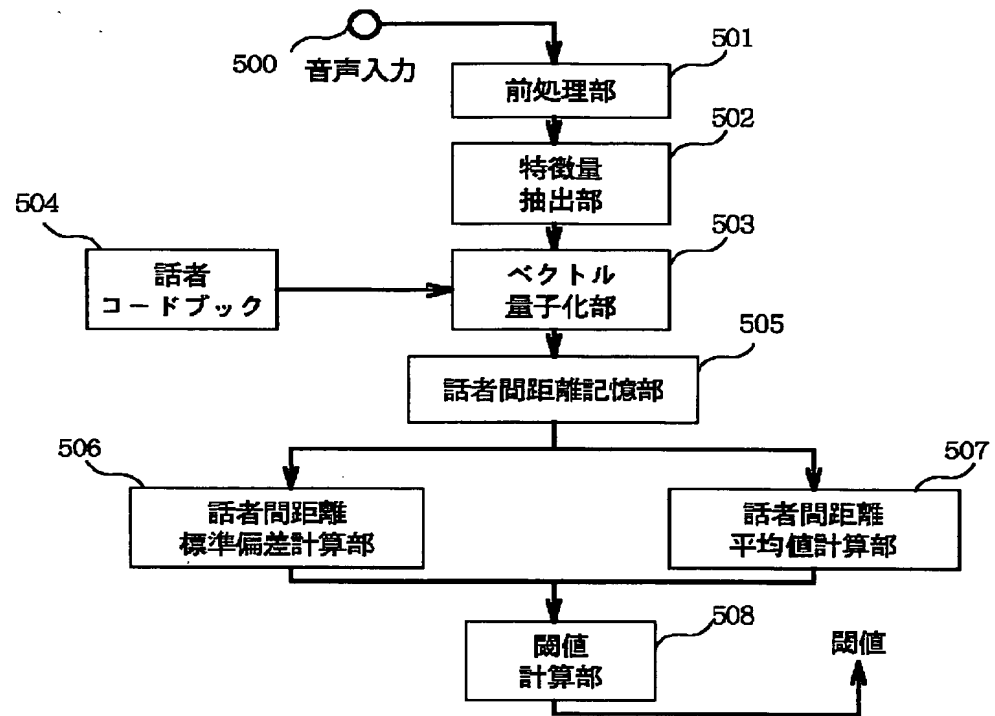


【図 4】

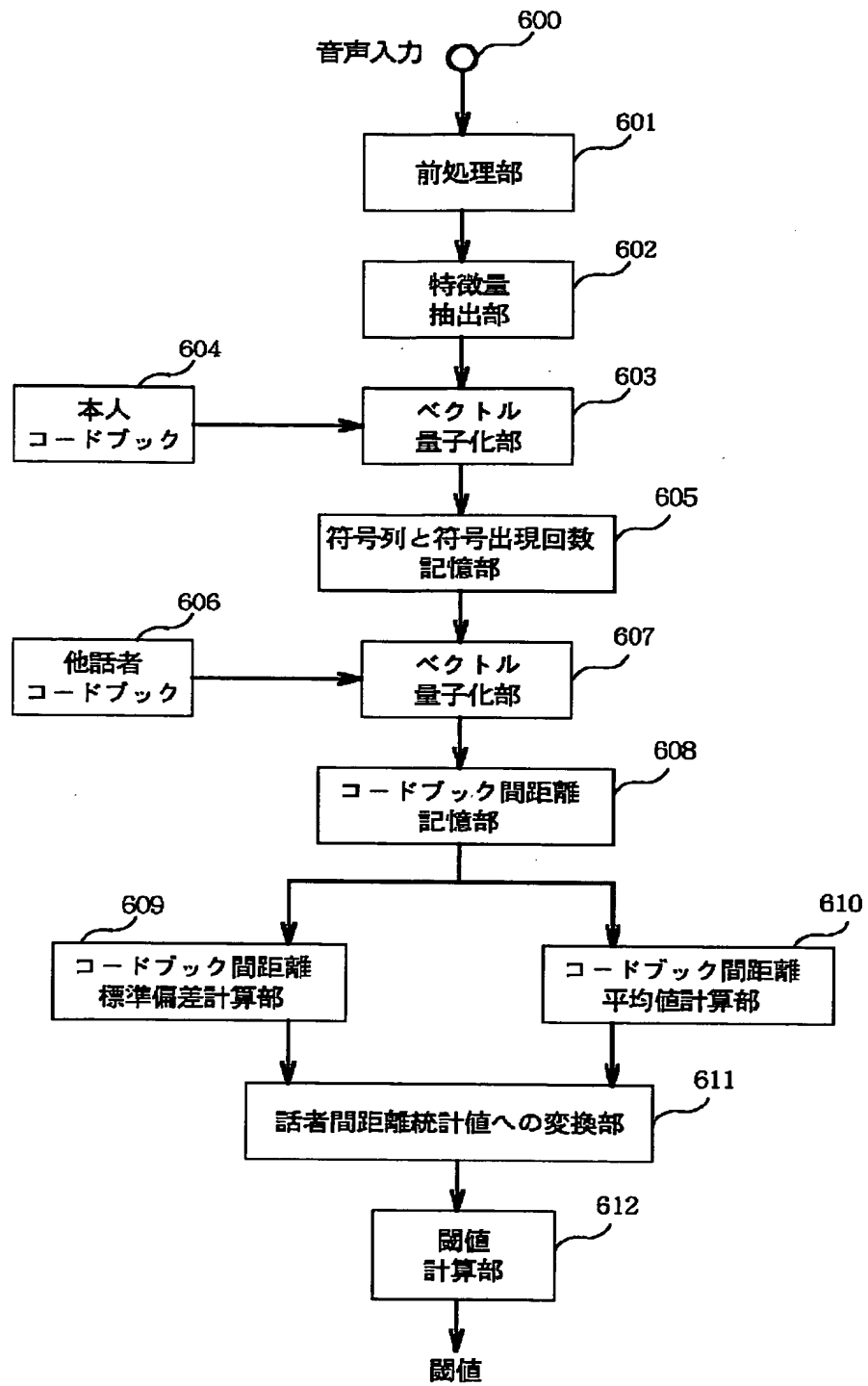




【図 5】



【図 6】



**Family list**

**2** family member for:

**JP4305699**

Derived from 1 application.

**1 METHOD AND DEVICE FOR RECOGNIZING VOICE**

Publication info: **JP3115016B2 B2** - 2000-12-04

**JP4305699 A** - 1992-10-28

---

Data supplied from the **esp@cenet** database - Worldwide

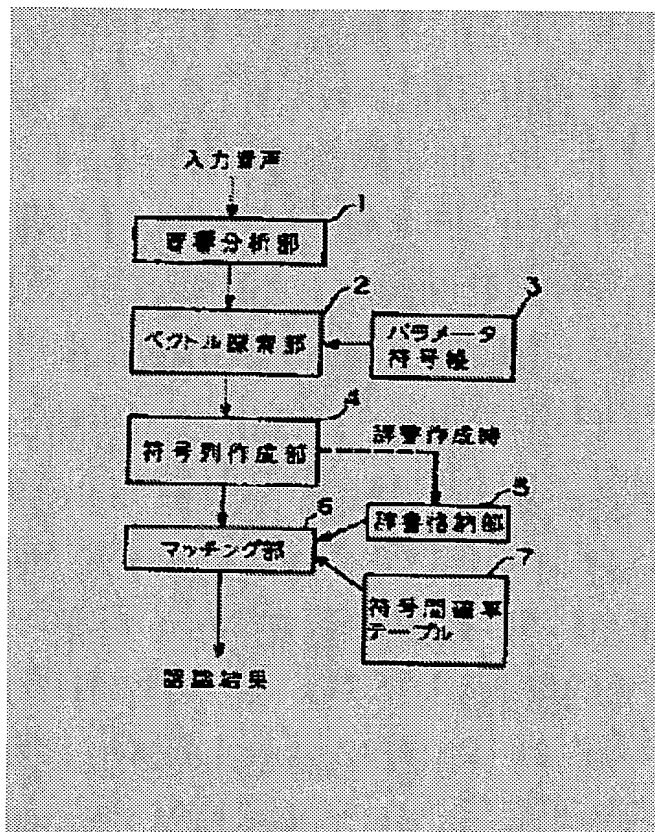
# METHOD AND DEVICE FOR RECOGNIZING VOICE

**Patent number:** JP4305699  
**Publication date:** 1992-10-28  
**Inventor:** MORII TOSHIYUKI; others: 02  
**Applicant:** MATSUSHITA ELECTRIC IND CO LTD  
**Classification:**  
- international: G10L3/00  
- european:  
**Application number:** JP19910071187 19910403  
**Priority number(s):**

## Abstract of JP4305699

**PURPOSE:** To recognize voices of many unspecified persons by using voices to be recognized that one and several speakers speak.

**CONSTITUTION:** In case of recognition, an acoustic analysis part 1 finds a feature vector from an input voice, frame by frame, a vector search part 2 calculates the distances to centroids stored in a parameter code book 3 to find the number of the centroid having the shortest distance, and a code string generation part 4 generates a code string. Then a matching part 6 finds the code string obtained from the input voice with the code string of a voice to be recognized which is stored in a dictionary storage part 5 and employs the voice having the highest similarity as a recognition result. The distance between codes at the time of the matching is found by looking up an inter-code probability table 7.



特開平4-305699

(43)公開日 平成4年(1992)10月28日

(51)Int.Cl.<sup>s</sup>

G 1 0 L 3/00

識別記号

3 0 1 D

庁内整理番号

8842-5H

F I

技術表示箇所

B 8842-5H

審査請求 未請求 請求項の数 3 (全 7 頁)

(21)出願番号 特願平3-71187

(22)出願日 平成3年(1991)4月3日

(71)出願人 000005821

松下電器産業株式会社

大阪府門真市大字門真1006番地

(72)発明者 森 井 利 幸

神奈川県川崎市多摩区東三田3丁目10番1号 松下技研株式会社内

(72)発明者 星 見 昌 克

神奈川県川崎市多摩区東三田3丁目10番1号 松下技研株式会社内

(72)発明者 二矢田 勝 行

神奈川県川崎市多摩区東三田3丁目10番1号 松下技研株式会社内

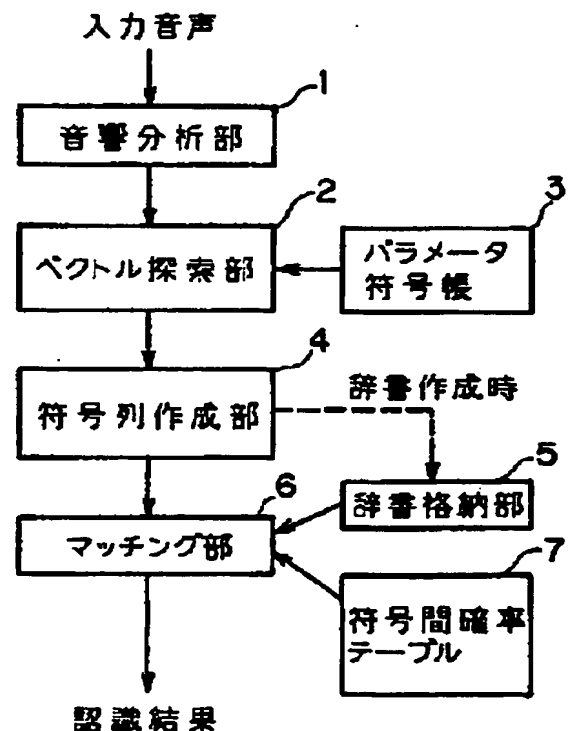
(74)代理人 弁理士 蔵合 正博

(54)【発明の名称】 音声認識方法および装置

(57)【要約】

【目的】 1名から数名の少数話者が発声した認識対象音声を用いて不特定多数の話者の音声を認識できるようにすること。

【構成】 認識時には、まず音響分析部1で入力音声からフレーム毎に特徴ベクトルを求め、次にベクトル探索部2でパラメータ符号帳3に格納されたセントロイドとの距離を計算し、最も距離の小さなセントロイドの番号を求めて符号列作成部4で符号列を作成する。次にマッチング部6で入力音声から得られた符号列と辞書格納部5に格納された認識対象音声の符号列とを照合し、最も類似度の大きい認識対象音声を認識結果とする。この照合の際の各符号間の距離は符号間確率テーブル7を参照することにより求める。



## 【特許請求の範囲】

【請求項1】 入力音声に対して音響分析を行ない、分析単位時間毎に得られる特徴パラメータの時系列である特徴ベクトルを求める音声分析手段と、多数の話者の音声を上記音声分析手段により分析することによって得られる特徴ベクトルの空間の代表ベクトルを格納するパラメータ符号帳と、入力音声から上記音声分析手段により得られる特徴ベクトルと上記パラメータ符号帳に格納された代表ベクトルとの距離を求め、最も近い代表ベクトルの符号を求めるベクトル探索手段とを備え、あらかじめ音素または音節位置のラベリングがなされている多数の話者の音声を上記音声分析手段とパラメータ符号帳とベクトル探索手段によって分析単位時間毎に符号化し、その符号とその符号を求めた分析単位時間に記された音素または音節のラベルとを用いて、各符号に符号化される特徴ベクトルが各音素または音節に属する確率を求めることにより確率列を作成しておき、入力音声の分析単位時間毎の特徴パラメータとして、入力音声から上記ベクトル探索手段によって得られた分析単位時間毎の符号に基づく上記確率列を用いて音声のマッチングを行なうことにより、不特定話者の音声の持つ個人性に影響されずに認識を行なう音声認識方法。

【請求項2】 入力音声に対して音響分析を行ない、分析単位時間毎に得られる特徴パラメータの時系列である特徴ベクトルを求める音声分析手段と、多数の話者の音声を上記音声分析手段により分析することによって得られる特徴ベクトルの空間の代表ベクトルを格納するパラメータ符号帳と、入力音声から上記音声分析手段により得られる特徴ベクトルと上記パラメータ符号帳に格納された代表ベクトルとの距離を求め、最も近い代表ベクトルの付号を求めるベクトル探索手段と、上記ベクトル探索手段により得られた符号を時間的に並べて符号列を作成する符号列作成手段と、1名から数名の話者が発声した認識対象音声を上記音声分析手段とパラメータ符号帳とベクトル探索手段と符号列作成手段により符号列に変換したものを基に作成された標準パターンを格納する辞書格納部と、上記パラメータ符号帳の符号間の類似度を格納する符号間確率テーブルと、上記入力音声を上記符号列作成手段により符号列に変換したものと上記辞書格納部に格納されている認識対象音声の標準パターンとしての符号列とを上記符号間確率テーブルに格納されている類似度を用いてマッチングして最も類似度の高い認識対象単語を認識結果とする音声マッチング手段とを有する音声認識装置。

【請求項3】 符号間確率テーブルが、あらかじめ音素または音節位置のラベリングがなされている多数の話者の音声を音声分析手段とパラメータ符号帳とベクトル探索手段によって分析単位時間毎に符号化し、その符号とその符号を求めた分析単位時間に記された音素または音節のラベルとを用いて、各符号に符号化される特徴ベク

トルが各音素または音節に属する確率を求め、さらにその確率の列を用いて符号化された2つの特徴ベクトルが同じ音素または音節に属する確率を求め、それらを格納することによって作成される請求項2記載の音声認識装置。

## 【発明の詳細な説明】

## 【0001】

【産業上の利用分野】本発明は、不特定話者が発声した単語音声を認識するための方法および装置に関する。

## 【0002】

【従来の技術】不特定話者を対象とした従来の音声認識技術については、たとえば、「ワードスポッティング手法を用いた不特定話者・少数語向け音声認識装置」(電子通信情報学会 SP88-18)に記載された方法が一般的である。

【0003】この方法では、入力された音声をまず音響分析し、音声の特徴パラメータの時系列に変換する。一方、あらかじめ認識装置側には、認識する単語毎にその単語の音声の特徴を示す単語標準パターンが用意されている。そして、話者の発声スピードを考慮して入力の特徴パラメータ列を時間的に伸縮しながら、ベイズ判定に基づく統計的距離尺度で単語標準パターンとのマッチングを行ない、最も距離が近いとされる単語を認識結果とする。この単語標準パターンは、その単語の平均的な特徴パラメータの時系列と、その分散行列によって構成されている。この分散行列によって話者の声の違いを吸収し、どんな話者の声でも認識することが可能となる。

【0004】上記の単語標準パターンの作成は以下の手順で行なう。

(1) 認識する単語集合(上記文献では10数字)について、330名が発声した単語音声データを収録し、音声データベースを作成する。

(2) 1つ1つの単語音声データに対して、スペクトル波形などのディスプレイ表示により、人間が目視で音声区間を検出し、単語の部分のみを切り出す。

(3) 切り出された区間を音響分析し、特徴パラメータ(LPC係数)の時系列を求め、さらに時間的間引きを行ない同じ時間長にする。これを単語パターンと呼ぶ。

(4) 得られた単語パターンを各単語毎に集め、各パラメータ列の平均と共分散行列とを求める。

【0005】このベイズ判定に基づく距離計算を行なうための標準パターンを作成するためには、上記のような多次元正規分布を仮定した統計分析が必要である。したがって、この構成の標準パターンは、数百名程度の多くの話者の音声を統計処理しなくては得られない。上記文献の例では、単語標準パターンを作成するために、330名の話者が発声した単語音声データを使用している。したがって、そのデータ作成には多大な労力が必要となる。

【0006】また、上記以外の不特定話者用音声認識の既存の方法としては、マルチ標準パターンを用いる方法が挙げられる。これは、1つの単語の標準パターンを代表的な単語パターン複数個により構成し、認識時には、この複数の単語パターンと入力パターンとの照合を行なうというものである。この方法は、複数のパターンを用いることによって不特定話者の音声認識しようとするものであるが、この複数のパターンを選択するためには、上記統計的距離尺度に基づく標準パターン作成時と同様に、多くの音声データと膨大な作業量とを必要とする。

#### 【0007】

【発明が解決しようとする課題】このように、既存の認識方法では、認識対象の音声の標準パターン作成に、音声データ収集や音声区間切り出し等のために多大な作業量が必要とする。したがって、認識対象の単語や文章を変更するのは大変困難であり、これは、語彙数が大きくなればなるほど深刻な問題となる。

【0008】本発明は、このような従来の問題を解決するものであり、1名から数名の少数話者が発声した認識対象音声を用いて不特定話者の音声認識を可能にするとともに、認識対象音声を容易に変更できる音声認識方法および装置を提供することを目的とする。

#### 【0009】

【課題を解決するための手段】本発明は、上記目的を達成するために、入力声を分析して得られる特徴パラメータの時系列である特徴ベクトルに対して、あらかじめ多数の話者で作成したパラメータ符号帳を用いて符号化を行ない、同様に符号列に変換された少数話者の音声パターンと符号間確率テーブルを用いてマッチングを行なうようにしたものである。

#### 【0010】

【作用】本発明は、上記構成により、まず入力声を分析して得られる特徴ベクトルに対して、多数の話者で作成したパラメータ符号帳を用いて符号化を行ない、分析単位時間（以下、フレームと呼ぶ。）毎に求めた符号を並べて符号列を作成する。そして、1名から数名の少数話者が発声した音声を同様に符号列に変換したものを基に作成した標準パターンとのマッチングを行ない、類似度を計算する。その際に用いられる符号間確率テーブルに格納された確率値は、多数の話者で作成した汎用性のある値であるので、個人性の影響を受けにくい。したがって、この確率値を基に単語の類似度を求めることによって、不特定話者の音声認識をすることができる。

【0011】また、どのような言葉も音素や音節の組合せで記述できるので、上記のパラメータ符号帳と符号間確率テーブルは1度作成しておけば十分であり、認識対象音声を変更しても常に同じものが使用できる。従って、不特定話者用の音声認識を行なうのに必要なものは、少数話者が発声した認識対象単語の音声データのみ

である。

【0012】以上により、簡単な手続で不特定話者用の音声認識が可能であり、かつ、語彙の変更に対して柔軟性のある認識装置の実現が可能になる。

#### 【0013】

【実施例】以下、本発明の実施例について説明するが、その前に本発明の基本的な考え方の背景について説明する。

【0014】人の声は有声音と無声音の2つに分類される。有声音は、声帯の振動として発せられた振動音が、調音器官と呼ばれる喉頭、咽頭、舌、あご、唇などで形成される声道を通る間に様々な変調を受けて、口から音声として出力されるという過程で発声される。すなわち、「あ」、「い」、「う」等の音韻性は声道の形状により与えられるのである。また、無声音は、音源が声帯でない場合もあるが、音韻性は有声音と同様に声道の形状によって決定される。しかし、声道を形成する喉、舌、歯、あご、唇等の形状や寸法は人毎に異なっているし、声帯の大きさも性別や年齢で異なる。このために、人毎に声の違いが生じることになる。つまり、人々による声の差異は調音器官の違いによるところが大きい。

【0015】一方、声が「あ」、「い」、「う」等の音韻としてでなく、単語や文として発せられるときは、声道の形は時間的に変化し、その変化によって言葉が形成される。たとえば、「赤い」(a k a i)と発声する場合、声道は、あごが開き舌の後方に決めるある/a/の発声から、喉頭部の閉鎖と急激な開放を伴う/k/に移り、更に再び/a/の形状に戻ってから徐々に舌を唇側に移動し、口を閉じた/i/に移る。このような声道の変化パターンは発声しようとしている言葉によって決まるものであり、人々による差異は少ないと考えられる。このように言葉としての音声を静的な声道の形状の違いとその時間的な変化に分離して考えると、前者は話者によってかなり異なるが、後者は比較的小さいと見ることが出来る。したがって、静的な声道の違いに基づく差異を何等かの方法で正規化できれば、不特定話者の音声認識が可能になる。

【0016】ところで、声道の形状の違いは、発せられた音声信号中では、周波数スペクトルの違いとして表現される。周波数スペクトルを話者間で正規化する最も単純な方法は、音素や音節などの短時間の音声標準パターンとのマッチングを行なって、発声された音声を音素や音節などの記号列にしてしまうことである。つまり、不特定話者用として作成された汎用の音素や音節の標準パターンを用いれば、話者の違いに大きく左右されずに、各音素や音節のどれに近いかという類似度情報を得ることができるのである。言換えると、周波数スペクトルをパターンマッチングによって音素や音節の類似度に変換することによって、話者の静的な声道の違いに基づく差

異を正規化することができるということである。そして、この正規化ができれば、声道の時間的変化は話者による差異が少ないのであるから、声道の変化パターンは、1人ないし数人分の音声データを上記正規化して得られる類似度の時間パターンにより作成することができる。したがって、少数話者の単語や文節の発声により、不特定話者用の音声標準パターンが得られる。

【0017】このような考え方にに基づき、本発明は次のように構成される。すなわち、予め多くの話者が発声した音声进行分析して、話者が発声する音全体の特徴パラメータの時系列である特徴ベクトルのセントロイド（重心）の集合を作成し、各セントロイドに番号を付けてパラメータ符号帳とする。また、その音声データに付加された音素位置のデータ（ラベルデータ）を利用して、上記パラメータ符号帳内の各セントロイドが各音素である確率を求めて、さらにその音素数の次元を持つ確率列から各々のセントロイドがお互いに同じ音素である確率を求めて、符号間確率テーブルを作成する。このテーブルに書かれた確率値は、話者の静的な声道の違いに基づく差異を受けにくい値である。標準パターンは、1名から数名の話者が発声した認識対象音声进行分析して得られる特徴ベクトルを上記パラメータ符号帳を用いて符号化し、セントロイドの番号の時系列（符号列）に変換することにより得られる。認識時には、入力音声に対して音響分析を行ない特徴ベクトルに変換した後、上記パラメータ符号帳によって符号化し符号列を求める。そして、標準パターンとしての符号列と照合を行なう。この際、符号間の距離は、上記符号間確率テーブルを参照するこ\*

$$d_j = \sum_{i=1}^{n_i} (x_i - a_i^{(j)})^2$$

$$\left\{ \begin{array}{l} \text{入力 Vector } x = (x_1, \dots, x_{n_i}) \\ \text{セントロイド Vector } a_j = (a_1^{(j)}, \dots, a_{n_i}^{(j)}) \end{array} \right.$$

【0022】図2はこのベクトル探索部2における探索の様子を示したものである。この図2の場合は、入力音声の特徴ベクトルに距離dが一番近い符号「1」に符号化される。そして、符号列作成部4において、各フレーム毎の番号を並べて符号列を作成する。

【0023】ここで、ベクトル探索部2において使用されるパラメータ符号帳3の作成方法について説明する。まず、多くの話者について、音韻バランスのとれた音声データを収録する。本実施例では多数の単語の音声データを使用している。次に、その音声の音声区間全てについて上記と同様の音響分析を行ない、各フレームの特徴ベクトルを求める。そして、それら全ての特徴ベクトルを集めて特徴ベクトルの母集団を作成し、さらに、この母集団に対してユークリッド距離に基づくクラスタリングを行ない、セントロイドを求めてパラメータ符号帳3を作成する。このクラスタリングは、母集団に対してサ

\*とにより求められる。

【0018】以下、本発明の一実施例について図面を参照して説明する。図1は本発明の一実施例の構成を示すものである。図1において、1は音響分析部、2はベクトル探索部、3はパラメータ符号帳、4は符号列作成部、5は辞書格納部、6はマッチング部、7は符号間確率テーブルである。

【0019】次に本実施例の動作について、最初に1名の話者の音声を辞書に登録する場合について説明する。図1において、まず入力音声に対して音響分析部1で1フレーム（本実施例では1フレーム=10msec）毎に線形予測分析（LPC分析）を行ない、特徴パラメータとしてLPCケプストラム係数（C0~C8 まで9個、C0は正規化残差パワー項で対数変換しておく。）を求める。

【0020】次に、ベクトル探索部2において、各フレームを中心とした特徴パラメータの時系列すなわち特徴ベクトルとパラメータ符号帳3に格納されている各セントロイド（重心）とのユークリッド距離の計算を以下の（数1）を用いて行ない、最も距離の近いセントロイドの番号を求める。特徴ベクトルは、本実施例では中心フレームから前4、後4フレームの計9フレーム分のLPCケプストラム係数（C0~C8）を1次元に並べたVector  $x = (C_0^{(1)}, C_1^{(1)}, \dots, C_8^{(1)}, C_0^{(2)}, C_1^{(2)}, \dots, C_8^{(2)}, \dots, C_0^{(9)}, C_1^{(9)}, \dots, C_8^{(9)})$  を意味する。

【0021】

【数1】

ンブルとセントロイド（重心）間のユークリッド距離が最小になるようなグループ分けを行ない、作成しようとする符号帳サイズの数のグループにわけて、そのグループのセントロイド（重心）で符号帳を作成する。

【0024】クラスタリングには幾つかの方法があり、本実施例に用いたクラスタリング・アルゴリズムは細胞分裂型のアルゴリズムである。このアルゴリズムを以下に順に示す。

(1)  $K=1$

(2) K個のグループのセントロイドを単純平均により求める。そして、それぞれのグループに属する全てのサンプルとセントロイドとのユークリッド距離を求め、その最大値をそのグループの歪とする。

(3) K個のグループの中で最も歪の大きいグループのセントロイドの附近に2つのセントロイドを作る（細胞分裂の核になる。）。



7

(4) K+1個のセントロイドを基にグループ分けを行ない、セントロイドを求め直す。

(5) 空のグループがあればそのセントロイドを抹消して(3)へ戻る。

(6) K+1個のグループの歪を(2)と同様に求め、その総和の変化量があらかじめ設定した微小なしきい値以下であれば(7)へ進み、しきい値より大きい場合は(4)へ戻る。

(7) K+1が目標のグループ数に達していなければK=K+1として(2)へ戻り、達していれば(8)へ進む。

(8) すべてのグループのセントロイドを求め、符号帳を作成する。

【0025】なお、上記アルゴリズムにおいて、本実施例におけるパラメータ符号帳3に格納されたセントロイドの数は全部で920個であり、収束検知に用いたしきい値は0.0001である。

【0026】ここで再び図1の実施例の説明に戻る。符号列作成部4において、各フレーム毎のセントロイドの番号を並べて符号列を作成した後、マッチング部6において、辞書格納部5に格納されている音声パターンとし\*

8

\*ての符号列とのマッチングを、符号間確率テーブル7とDPマッチングを用いて行ない、各音声の類似度を求める。そして、各類似度を比較し、最も高いものを認識結果として出力する。この辞書格納部5、マッチング部6、符号間確率テーブル7について、以下に説明する。

【0027】まず、辞書格納部5に格納される音声パターンの作成手順を述べる。最初に認識対象音声について1人の話者の音声を収録する。次に、認識時と同様に音響分析を行ない、特徴ベクトルを求める。さらに、認識時と同様にパラメータ符号帳を用いて符号化を行ない、各音声の符号列を求める。そして、この符号列を音声のパターンとして辞書格納部5に格納する。

【0028】次に、マッチング部6について述べる。入力音声1から得られる符号列と辞書格納部5に格納されている音声パターンは、一般にその長さが異なっている。そこで、このマッチングをDPマッチングを用いて行なう。本実施例で用いた漸化式の例を(数2)に示す。

【0029】

【数2】

$$g(i, j) = \max \begin{bmatrix} g(i-2, j-1) + 3L(i, j) \\ g(i-1, j-1) + 2L(i, j) \\ g(i-1, j-2) + 3L(i, j) \end{bmatrix}$$

【0030】ここで、辞書側のフレーム番号がj、入力のフレーム番号がi、第iフレームと第jフレームの類似度がL(i, j)、累積類似度がg(i, j)である。類似度L(i, j)は、辞書側のj番目にある符号と、入力のi番目の符号(セントロイドの番号)を基に、符号間確率テーブル7を参照して求める。

【0031】符号間確率テーブル7は、図3に示すように、各符号間の類似度(同じである確率で、図3には生の値を記しているが、実際には対数をとってある。)が入っており、マッチングの際には、比較する符号を縦横に見てその間の確率値を類似度として用いる。この値は、2つのフレームが同じである確率であり、このDPマッチングの結果得られる累積類似度は、マッチングパスにおいて対応する全てのフレームが同じである確率になることに注意すべきである。この「2つのフレームが同じである確率」については、以下の符号間確率テーブル7についての説明の中でその意味を述べる。

【0032】次に、符号間確率テーブル7について、その考え方と作成法について説明する。ベクトル探索部2において求められる符号は、特徴ベクトルのベクトル空間(ユークリッド空間)上における大まかな位置を示すものである。また、上記課題を解決するための手段の項で述べたように、少数の話者の音声データを不特定話者の標準パターンに変換するためには、多数話者の音声

データから作成した音素や音節の標準パターンとマッチングを行なう必要がある。そこで、本実施例では音素を基本単位とした統計分析により、その符号に符号化される特徴ベクトルが各音素に属する確率を求める。

【0033】まず、音素の位置(始端と終端)がラベル付けされている音声データを上記認識時と同様に音響分析して、各フレームの特徴ベクトルを求めた後、上記パラメータ符号帳3によって符号化し(最も近いセントロイドの番号を求める。)、各フレームの符号(セントロイドの番号)を求める。次に、ラベルを参照することにより、そのフレームが何の音素に属しているかがわかるので、各セントロイド毎にその音素数分のエリアを用意し、そのセントロイド番号になった特徴ベクトルの音素番号のエリアに加算していく。その結果、図4のように、各セントロイド番号のエリアには、その符号になった特徴ベクトルが各音素であった個数が入っている。たとえば、図4の符号「3」を例としてみると、多数話者の全ての音声から得られた多くの特徴ベクトルのうち、「3」に符号化されたものは全部で1200個あり、このうち31個が/a/の音素であり、また、40個が/o/、935個が/u/であったことを示している。そこで、この個数を全体で割れば、その番号に符号化された特徴ベクトルが各音素になる確率が得られる。これを確率列と呼ぶ。例えば、この例では、「3」に符号化さ

9

れた特徴ベクトルが/a/である確率は0.0258であり、/u/である確率は0.779である。この値はいずれも多数話者から求めた汎用性のある確率であるから、音声はこの確率列に変換することによって、話者の声道の違いに基づく差異を正規化できる。したがって、入力音声の特徴パラメータとしてこの確率列を用いることは、個人性に影響されにくいという点で大変有効である。

【0034】そして、さらにこの確率列を用いて2つの符号の類似度を求める。符号iが音素jに属する確率をPijとすると、符号mと符号nが同じ音素である確率D<sub>mn</sub>は次の(数3)によって求めることができる。

【0035】

【数3】

$$D_{mn} = \sum_{j=1}^J P_{mj} \cdot P_{nj}$$

【0036】最後に、このD<sub>mn</sub>をマトリックスに表現して図3のような符号間確率テーブルを作成する。このD<sub>mn</sub>は、人の違いによらず、その符号と符号が音素としてどれだけ似ているかという程度をあらわす値である。したがって、上記した辞書格納部5の説明文中における「同じである確率」とは、この「同じ音素である確率」に相当する。

【0037】なお、本実施例における音素とは、/a/, /o/, /u/, /i/, /e/, /j/, /w/, /m/, /n/, /

【0038】

【外1】

0

【0039】/ (語中), /b/, /d/, /g/ (語頭), /r/, /z/, /h/, /s/, /c/, /p/, /t/, /k/, Q (促音), /=/ (撥音) の23音素とした。

【0040】以上が、本実施例における1人の話者の音声パターンを標準パターンとする音声認識方法である。次に、複数の話者で標準パターンをつくる方法は2つある。1つは、一人一人の音声パターンをそのままマルチパターンとしてマッチングに用いる方法で、この場合は認識時の計算量はモデルとする話者数に比例して増加するが、より話者に適応した認識を行なうことができ、認識性能を向上させることができる。2つ目は、それぞれの音声パターンの長さをDPマッチングを用いて正規化してから平均化する方法である。この場合、平均化する段階は、特徴ベクトルの段階と確率列の段階の2種類があり、どちらにしても、より安定した標準パターンを得ることができ、認識性能を向上させることができる。

【0041】ここで、本発明の有効性を検証するために、本実施例を用いて単語認識実験を行なった。単語数は212個である。音素位置のラベル付けがなされた20人(男女各10人)の話者の音声データを用いて、バ

10

ラメータ符号帳と符号間確率テーブルを作成し、この内の1人(男性話者)の212単語の単語音声を用いて音声の標準パターンを作成した。認識対象は上記話者以外の話者20人(男女各10人)の212単語音声である。実験の結果、平均90.1%(男性89.39%、女性90.81%)という高い単語認識率が得られた。男性の音声パターンを用いているが、男女の認識率の差はほとんど無い。したがって、ベクトル探索と符号間確率による距離計算により、声の個人性が吸収されており、1人の話者でも不特定用の標準パターンが得られることが検証できた。

【0042】このように、本実施例では、入力音声を音響分析することによって得られる特徴パラメータの時系列である特徴ベクトルを用い、そのままその空間上でマッチングするのではなく、その特徴ベクトルを多数話者で作成したパラメータ符号帳で符号化し、マッチングの際に符号間確率テーブルを参照して符号間類似度を求めることにより、少数話者が発声した音声を登録するだけで不特定話者の音声を精度良く認識することができるようになる。

【0043】

【発明の効果】以上のように、本発明は、入力声を分析して得られる特徴パラメータの時系列である特徴ベクトルに対して、あらかじめ多数の話者で作成したパラメータ符号帳を用いて符号化を行ない、同様に符号列に変換された少数話者の音声パターンと符号間確率テーブルを用いてマッチングを行なうことによって、1人から数名の少数話者が発声した認識対象単語を辞書として登録するだけで辞書が更新でき、また、高い音声認識率を得ることができる。

【0044】このように、本発明は、不特定話者用音声認識装置の性能向上および様々な用途に適用するために、標準パターンを作成するための労力削減に対して極めて大きく貢献することができる。

【図面の簡単な説明】

【図1】本発明の一実施例における音声認識装置の構成を示すブロック図

【図2】同実施例におけるベクトル探索部の機能を説明するための模式図

【図3】同実施例における符号間確率テーブルの一例を示す説明図

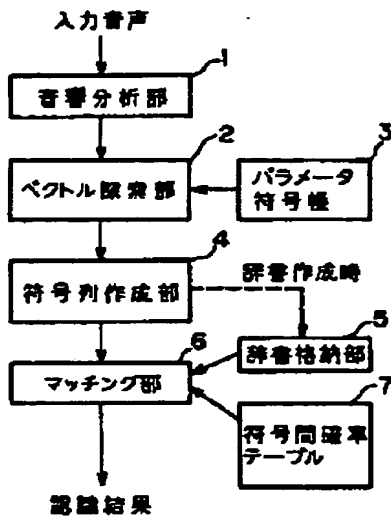
【図4】同実施例における各セントロイドが各音素であった個数を示す説明図

【符号の説明】

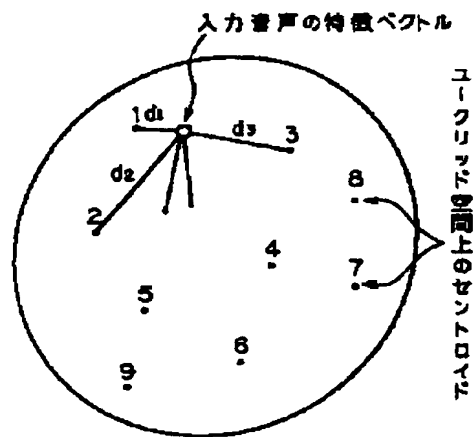
- 1 音響分析部
- 2 ベクトル探索部
- 3 パラメータ符号帳
- 4 符号列作成部
- 5 辞書格納部
- 6 マッチング部

## 7 符号間確率テーブル

【図1】



【図2】



【図4】

【図3】

符号(セントロイド番号)

	1	2	3	4	5	6	7	...
1	0.8	0.02	0.03	0.012	0.001	0.01	0	
2	0.02	0.5	0	0.11	0.001	0.001	0.01	
3	0.03	0	0.7	0.05	0.001	0.001	0.001	
4	0.012	0.11	0.05	0.6	0	0.03	0.03	
5	0.001	0.001	0.001	0	0.7	0.01	0.02	
...								

(単語フェィローレンゼ) 単語

符号(セントロイド番号)

各音素であった回数

	1	2	3	4	5	6	...
/a/	131	363	31	0	13	15	
/o/	22	56	40	10	94	3	
/u/	5	81	935	4	106	0	
/i/	13	29	44	3	40	0	
/e/	8	16	3	2	35	9	
.							
.							
.							
計	260	600	1200	120	424	51	

**Family list**

**1** family member for:

**JP10091186**

Derived from 1 application.

**1 VOICE RECOGNIZING METHOD**

Publication Info: **JP10091186 A** - 1998-04-10

---

Data supplied from the **esp@cenet** database - Worldwide

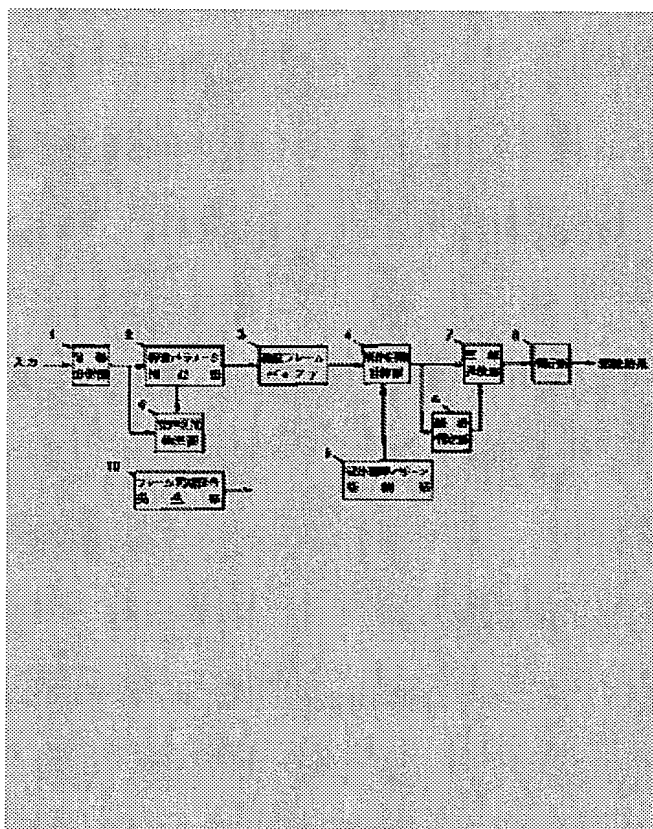
# VOICE RECOGNIZING METHOD

**Patent number:** JP10091186  
**Publication date:** 1998-04-10  
**Inventor:** FUTAYADA KATSUYUKI; HOSHIMI MASAKATSU; HIRAOKA SEIJI; KIMURA TATSUYA  
**Applicant:** MATSUSHITA ELECTRIC IND CO LTD  
**Classification:**  
- **international:** G10L3/00  
- **europaen:**  
**Application number:** JP19970295111 19971028  
**Priority number(s):**

## Abstract of JP10091186

**PROBLEM TO BE SOLVED:** To obtain a voice recognizing method which is tolerative for increasing of the number of vocabulary and noise and a recognition rate is high by accumulating distance obtained from an input vector and a partial pattern with a statistical distance scale, obtaining the accumulated distance, and making a word of the minimum accumulation distance the recognized result from the accumulated distance.

**SOLUTION:** A voice recognizing device is constituted with an acoustic analyzing section 1, a feature parameter extracting section 2, a plural frame buffer 3, a partial distance calculating section 4, a partial standard pattern storing section 5, a path discriminating section 6, a distance accumulating section 7, a discriminating section 8, a voice section detecting section 9, and a frame synchronizing signal generation section 10. And a partial distance between an input vector formed by plural frames and a partial (standard) pattern of a voice of a word is obtained with a statistical distance scale based on posterior probability, the input vector is updated shifting a frame, distance between each partial vector is accumulated, and a word in which



accumulated distance is minimum is made the recognized result. Thereby, a high recognition rate can be obtained for voice recognition for an unspecified talker.

---

Data supplied from the **esp@cenet** database - Worldwide



## 【特許請求の範囲】

【請求項1】 認識対象単語の標準パターンを部分パターンの接続で作成する工程と、入力音声からフレームをシフトしながら入力ベクトルを求める工程と、前記入力ベクトルと前記部分パターンとから統計的距離尺度で求めた距離を累積し累積距離を求める工程と、前記累積距離から最小累積距離の単語を認識結果とする工程とを有することを特徴とする音声認識方法。

【請求項2】 多数の人が発声した音声データを用いて、認識対象単語を部分区間に分割し、その部分区間を表現する部分（標準）パターンを接続した認識対象単語の標準パターンを、全ての認識対象単語に対して予め生成する工程と、入力音声を一定時間長（フレーム）ごとに分析して特徴パラメータを求め、複数フレームの特徴パラメータで入力ベクトルを求める工程と、前記入力ベクトルと前記各部分パターンとの部分距離を統計的距離尺度で求める工程と、フレームをシフトしながら生成した入力ベクトルと前記部分パターンとの部分距離を累積した累積距離を求める工程と、全認識対象単語の標準パターンに対する累積距離を相互に比較して、最小累積距離の単語を認識結果とする工程とを有することを特徴とする音声認識方法。

【請求項3】 認識対象単語の部分パターンは、互いに重なる区間（フレーム）を含むことを特徴とする請求項1または2記載の音声認識方法。

【請求項4】 多数の人が発声した音声データを用いて、認識対象単語を複数フレームからなる部分区間に分割し、その部分区間を表現する部分（標準）パターンを接続した認識対象単語の標準パターンを、全ての認識対象単語に対して予め生成する工程と、入力音声を一定時間長（フレーム）ごとに分析して特徴パラメータを求め、複数フレームの特徴パラメータで入力ベクトルを求める工程と、前記入力ベクトルと前記各部分パターンとの部分距離を事後確率に基づく統計的距離尺度で求める工程と、フレームをシフトしながら生成した入力ベクトルと前記部分パターンとの部分距離を累積した累積距離を求める工程と、全認識対象単語の標準パターンに対する累積距離を相互に比較して、最小累積距離の単語を認識結果とする工程とを有することを特徴とする音声認識方法。

## 【発明の詳細な説明】

## 【0001】

【発明の属する技術分野】本発明は人間の声を機械に認識させる音声認識方法に関するものである。

## 【0002】

【従来の技術】近年、使用者の声を登録することなしに、誰の声でも認識できる不特定話者用の認識装置が実用として使われるようになった。不特定話者用の実用的な方法として、本出願人が、以前に出願した2つの特許（特開昭61-188599号公報、特開昭62-111293号公報）を

従来例として説明する。特開昭61-188599号公報を第1の従来例、特開昭62-111293号公報を第2の従来例とする。

【0003】第1の従来例の方法は入力音声の始端、終端を求めて音声区間を決定し、音声区間を一定時間長に（Iフレーム）に線形伸縮し、これと単語標準パターンとの類似度を統計的距離尺度を用いてパターンマッチングをすることによって求め、単語を認識する方法である。

10 【0004】単語標準パターンは、認識対象単語を多くの人に発声させて音声サンプルを収集し、すべての音声サンプルを一定時間長Iフレーム（実施例ではI=16）に伸縮し、その後、単語ごとに音声サンプル間の統計量（平均値ベクトルと共分散行列）を求め、これを加工することによって作成している。すなわち、すべての単語標準パターンの時間長は一定（Iフレーム）であり、原則として1単語に対し1標準パターンを用意している。

20 【0005】第1の従来例では、パターンマッチングの前に音声区間を検出する必要があるが、第2の従来例は音声区間検出を必要としない部分が異なっている。パターンマッチングによって、ノイズを含む信号の中から音声の部分抽出して認識する方法（ワードスポッティング法）を可能とする方法である。すなわち、音声を含む十分長い入力区間内において、入力区間内に部分領域を設定し、部分領域を伸縮しながら標準パターンとのマッチングを行なう。そして、部分領域を入力区間内で単位時間ずつシフトして、また同様に標準パターンとのマッチングを行なうという操作を設定した入力区間内全域で

30 行ない、すべてのマッチング計算において距離が最小となった単語標準パターン名を認識結果とする。ワードスポッティング法を可能にするために、パターンマッチングの距離尺度として事後確率に基づく統計的距離尺度を用いている。

## 【0006】

【発明が解決しようとする課題】従来例の方法は、小型化が可能な実用的な方法であり、特に第2の従来例は、騒音にも強いことから実用として使われ始めている。

40 【0007】しかし、従来技術の課題は、十分な単語認識率が得られないことである。このため、語彙の数が少ない用途にならば使うことが出来るが、語彙の数を増やすと認識率が低下して実用にならなくなってしまう。従って、従来技術の方法では認識装置の用途が限定されてしまうという課題があった。

【0008】本発明は上記従来の課題を解決するもので、語彙数の増加や騒音に対して頑強な認識率の高い音声認識方法を提供することを目的とするものである。

## 【0009】

50 【課題を解決するための手段】この課題を解決するために本発明は、多数の人が発声した音声データを用いて、



認識対象単語を隣接するフレームを共有する部分区間に分割し、その部分区間を表現する部分（標準）パターンを接続した認識対象単語の標準パターンを、全ての認識対象単語に対して予め生成する工程と、入力音声を一定時間長（フレーム）ごとに分析して特徴パラメータを求め、複数フレームの特徴パラメータで入力ベクトルを求める工程と、前記入力ベクトルと前記各部分パターンとの部分距離を事後確率に基づく統計的距離尺度で求める工程と、フレームをシフトしながら生成した入力ベクトルと前記部分パターンとの部分距離を累積した累積距離を求める工程と、全認識対象単語の標準パターンに対する累積距離を相互に比較して、最小累積距離の単語を認識結果とする工程とを有するものである。

【0010】このことにより、語彙数の増加や騒音に対して頑強で認識率の高い音声認識方法が得られる。

【0011】

【発明の実施の形態】本発明の請求項1に記載の発明は、認識対象単語の標準パターンを部分パターンの接続で作成する工程と、入力音声からフレームをシフトしながら入力ベクトルを求める工程と、前記入力ベクトルと前記部分パターンとから統計的距離尺度で求めた距離を累積し累積距離を求める工程と、前記累積距離から最小累積距離の単語を認識結果とする工程とを有するもので、フレームをシフトしながら入力音声から求めた入力ベクトルと、単語音声の標準パターンを構成する部分（標準）パターンとの部分距離を統計的距離尺度で求め、その距離を累積し、最小累積距離の単語を認識結果とするもので、不特定話者用の音声認識に対して認識率が得られものである。

【0012】請求項2の発明は、多数の人が発声した音声データを用いて、認識対象単語を部分区間に分割し、その部分区間を表現する部分（標準）パターンを接続して認識対象単語の標準パターンを、全ての認識対象単語に対して予め生成する工程と、入力音声を一定時間長（フレーム）ごとに分析して特徴パラメータを求め、複数フレームの特徴パラメータで入力ベクトルを求める工程と、前記入力ベクトルと前記標準パターンを構成する各部分パターンとの部分距離を統計的距離尺度で求める工程と、フレームをシフトしながら生成した入力ベクトルと前記部分パターンとの部分距離を累積した累積距離を求める工程と、全認識対象単語の標準パターンに対する累積距離を相互に比較して最小累積距離の単語を認識結果とする工程とを有するもので、複数のフレームで形成される入力ベクトルと、単語音声を部分区間に分割し、その部分区間を表現する部分（標準）パターンとの部分距離を事後確率に基づく統計的距離尺度で求め、フレームをシフトしながら入力ベクトルを更新して各部分ベクトルととの間の距離を累積し、累積距離を最小とする単語を認識結果とするもので、不特定話者用の音声認識において、語彙数の増加や騒音に対して頑強で高い認識率が

得られ、また処理が単純なので、信号処理プロセッサ（DSP）等を用いて、小型でリアルタイム動作が可能な認識装置を実現するという作用を有する。

【0013】請求項3記載の発明は、請求項1または2において、認識対象単語の部分区間は、互いに重なる区間を含むように分割するもので、区間の境界の動き情報を確実に得ることができ、より詳細な部分パターンが生成できるという作用を有する。

【0014】請求項4記載の発明は、多数の人が発声した音声データを用いて、認識対象単語を複数フレームからなる部分区間に分割し、その部分区間を表現する部分（標準）パターンを接続した認識対象単語の標準パターンを、全ての認識対象単語に対して予め生成する工程と、入力音声を一定時間長（フレーム）ごとに分析して特徴パラメータを求め、複数フレームの特徴パラメータで入力ベクトルを求める工程と、前記入力ベクトルと前記各部分パターンとの部分距離を事後確率に基づく統計的距離尺度で求める工程と、フレームをシフトしながら生成した入力ベクトルと前記部分パターンとの部分距離を累積した累積距離を求める工程と、全認識対象単語の標準パターンに対する累積距離を相互に比較して、最小累積距離の単語を認識結果とする工程とを有するもので、複数フレームからなる部分パターンとし、入力ベクトルと部分距離を求める際に事後確率に基づく統計的距離尺度で求めることにより、入力的位置や部分パターンの違いにもかかわらず部分距離を求めることができるという作用を有する。

【0015】以下、本発明の実施の形態について、図面を用いて説明する。実施の形態1は、入力音声の始端、終端があらかじめ検出されている場合における実施例である。この場合は音声区間でのみパターンマッチングを行えばよい。また、実施の形態2は、入力音声の始端、終端が未知の場合の実施例である。この場合は入力音声を含む十分広い区間内を対象として、入力信号と標準パターンのマッチングを区間全域にわたって単位時間ずつシフトしながら行ない、距離が最小となる部分区間を切り出す方法を用いる。この種の方法を一般的にワードスポッティングと呼んでいる。

【0016】（実施の形態1）図1に、本発明の実施の形態1の音声認識装置の機能ブロック図を示し、説明する。

【0017】図1において、音響分析部1は入力信号をAD変換して取込み（サンプリング周波数10kHz）、一定時間長（フレームと呼ぶ。本実施例では10ms）ごとに分析する。本実施例では線形予測分析（LPC分析）を用いる。特徴パラメータ抽出部2では分析結果に基づいて、特徴パラメータを抽出する。本実施例では、LPCケプストラム係数（ $C_0 \sim C_{10}$ ）および差分パワー値 $V_0$ の12個のパラメータを用いている。入力の1フレームあたりの特徴パラメータを

5

【0018】

【外1】

 $x_j$ 

【0019】と表すことにすると、特徴パラメータは(数1)のようになる。

【0020】

【数1】

$$x_j = (V_0, C_0, C_1, \dots, C_p)$$

【0021】ただし、jは入力フレーム番号、pはケプストラム係数の次数である(p=10)。フレーム同期信号発生部10は、10msごとに同期信号を発生する部分であり、その出力は全てのブロックに入る。即ち、システム全体がフレーム同期信号に同期して作動する。

【0022】音声区間検出部9は、入力信号音声の始

$$x_j = (x_{j-L_1}, x_{j-L_1+m}, x_{j-L_1+2m}, \dots, x_j, x_{j+m}, \dots, x_{j+L_2})$$

【0027】すなわち、上記入力ベクトルはmフレームおきにj-L<sub>1</sub>~j+L<sub>2</sub>フレームの特徴パラメータを統合したベクトルである。L<sub>1</sub>=L<sub>2</sub>=3, m=1 とすると上記入力ベクトルの次元数は (P+2) × (L<sub>1</sub>+L<sub>2</sub>+1) = 12×7=84となる。なお、(数2)ではフレーム間隔mは一定になっているが、必ずしも一定である必要はない。mが可変の場合は非線形にフレームを間引くことに相当する。

【0028】部分標準パターン格納部5は、認識対象とする各単語の標準パターンを、部分パターンの結合として格納してある部分である。ここで、本実施例における標準パターン作成法を、やや詳細に説明する。

【0029】話をわかり易くするために、今、認識対象単語を日本語の数字「イチ」「ニ」「サン」「ヨン」「ゴ」「ロク」「ナナ」「ハチ」「キュウ」「ゼロ」の10種とする。このような例を用いても説明の一般性にはなんら影響はない。

【0030】たとえば、「サン」の標準パターンは次のような手順で作成する。

(1) 多数の人(100名とする)が「サン」と発声したデータを用意する。

【0031】(2) 100名の「サン」の持続時間分布を調べ、100名の平均時間長I<sub>3</sub>を求める。

【0032】(3) 時間長のI<sub>3</sub>サンプルを100名の中から探し出す。複数のサンプルがあった場合はフレームごとに複数サンプルの平均値を計算する。このように求められた代表サンプルを(数3)で示す。

【0033】

【数3】

$$S_0 = (s_1, s_2, \dots, s_i, \dots, s_{I_3})$$

【0034】ここで

※

$$x_{(i)}^n = (x_{(i)-3}^n, x_{(i)-2}^n, \dots, x_{(i)}^n, x_{(i)+1}^n, \dots, x_{(i)+3}^n)$$

【0043】ここで、(i)は第n番目のサンプル中、代

(4)

6

\*端、終端を検出する部分である。音声区間の検出法は音声のパワーを用いる方法が簡単で一般的であるが、どのような方法でもよい。本実施例では音声の始端が検出された時点で認識が始まり、j=1になる。

【0023】複数フレームバッファ3は、第jフレームの近隣のフレームの特徴パラメータを統合して、パターンマッチング(部分マッチング)に用いる入力ベクトルを形成する部分である。すなわち、第jフレームに相当する入力ベクトル

【0024】

【外2】

 $x_j$ 

【0025】は、次式で表わされる。

【0026】

【数2】

※【0035】

【外3】

s

20

【0036】は1フレームあたりのパラメータベクトルであり、(数1)と同様に11個のLPCケプストラム係数と差分パワーで構成される。

【0037】(4) 100名分のサンプルの1つ1つと代表サンプルとの間でパターンマッチングを行ない、代表サンプルと100名分の各サンプルとの間の対応関係(最も類似したフレーム同士の対応)を求める。距離計算はユークリッド距離を用いる。代表サンプルのiフレームと、あるサンプルのi'フレームとの距離d<sub>i, i'</sub>は(数4)で表わされる。

【0038】

【数4】

$$d_{i, i'} = (x_{i'} - s_i)^t \cdot (x_{i'} - s_i)$$

【0039】ここで、tは転置行列であることを表す。

なお、フレーム間の対応関係はダイナミックプログラミング(DP法)の手法を用いれば効率よく求めることができる。

【0040】(5) 代表サンプルの各フレーム(i=1~I<sub>3</sub>)に対応して、100名分のサンプルそれぞれから(数2)の形の部分ベクトルを切出す。簡単化のためL<sub>1</sub>=L<sub>2</sub>=3, m=1 とする。

【0041】代表サンプルの第iフレームに相当する、100名の中の第n番目のサンプルの部分ベクトルは以下ようになる。

【0042】

【数5】

7

とを示す。

【0044】

【外4】

 $x_i$ 

【0045】は本実施例では84次元のベクトルである  
( $n=1\sim 100$ )。

(6) 100名分の上記ベクトルの平均値

【0046】

【外5】

 $\mu_k^i$ 

【0047】(本例では $k=3$ ; 84次元)と共分散行列

【0048】

【外6】

 $w_k^i$ 

【0049】(84×84次元)を求める( $i=1\sim I_3$ )。平均値と共分散行列は標準フレーム長の数 $I_3$ だけ存在することになる(ただし、これらは必ずしも全フレームに対して作成する必要はない。間引いて作成してもよい。)

\*20

$$W = \frac{1}{I_1 + I_2 + \dots + I_k + g} \left( \sum_{k=1}^K \sum_{i=1}^{I_k} W_k^i + g \cdot W_e \right)$$

【0058】ここで $K$ は認識対象単語の種類( $K=10$ )、 $I_k$ は単語 $k$ ( $k=1, 2, \dots, K$ )の標準時間長を表す。また、 $g$ は周囲パターンを混入する割合であり通常 $g=1$ とする。

【0059】b. 各単語の部分パターン

【0060】

【外9】

 $A_k^i$ 

【0061】及び

【0062】

※

$$B_k^i = \mu_k^i \cdot W^{-1} \cdot \mu_k^i - \mu_e \cdot W^{-1} \cdot \mu_e$$

【0066】これらの式の導出は後述する。図2に標準パターン作成法の概念図を示す。図2(a)は入力信号が「サン」の場合の音声のパワーパターンを示す。図2(b)は部分パターンの作成法を概念的に示したものである。音声サンプルの始端と終端の間において、代表サンプルとのフレーム対応を求めて、それによって音声サンプルを $I_3$ に分割する。図では代表サンプルとの対応フレームを(i)で示してある。そして、音声の始端(i)=1から終端(i)= $I_3$ の各々について、(i)- $L_1 \sim (i)+L_2$ の区間の100名分のデータを用いて平均値と共分散を計算し、部分パターン

【0067】

【外11】

 $A_k^i$ 

【0068】

8

\*【0050】上記(1)～(6)と同様の手続きで「サン」以外の単語に対しても84次元のベクトルと共分散行列を求める。

【0051】そして、全ての単語に対する100名分すべてのサンプルデータに対し、移動平均

【0052】

【外7】

 $\mu_e$ 

【0053】(84次元)と移動共分散行列

【0054】

【外8】

 $w_e$ 

【0055】(84×84次元)を求める。これらを周囲パターンと呼ぶ。次に平均値と共分散を用いて標準パターンを作成する。

【0056】a. (数6)により共分散行列を共通化する。

【0057】

【数6】

※【外10】

 $B_k^i$ 

【0063】を作成する。

【0064】

【数7】

$$A_k^i = 2 \cdot W^{-1} \cdot (\mu_k^i - \mu_e)$$

30

【0065】

【数8】

【外12】

 $B_k^i$ 

【0069】を求める。従って、単語 $k$ の標準パターンは互にオーバーラップする区間を含む $I_k$ 個の部分パターンを接続して(寄せ集めた)ものになる。図2(c)は周囲パターンの作成方法を示す。周囲パターンは標準パターン作成に使用した全データに対して、図のように $L_1+L_2+1$ フレームの部分区間を1フレームずつシフトさせながら移動平均値と移動共分散を求める。周囲パターン作成の範囲は音声区間内のみならず、前後のノイズ区間も対象としてもよい。後述する第2の実施例では周囲パターンにノイズ区間を含める必要がある。

【0070】次に部分距離の計算について述べる。上記のようにしてあらかじめ作成されている各単語の部分標

準パターンと複数フレームバッファ3との間の距離（部分距離）を部分距離計算部4において計算する。

【0071】部分距離の計算は、(数2)で示す複数フレームの情報を含む入力ベクトルと各単語の部分パターンとの間で、統計的な距離尺度を用いて計算する。単語全体としての距離は部分パターンとの距離（部分距離と呼ぶ）を累積して求めることになるので、入力的位置や部分パターンの違いにかかわらず、距離値が相互に比較できる方法で部分距離を計算する必要がある。このためには、事後確率に基づく距離尺度を用いる必要がある。

(数2)の形式の入力ベクトルを

$$P(\omega_k | X) = P(\omega_k) \cdot P(X | \omega_k) / P(X)$$

【0077】右辺第1項は、各単語の出現確率を同じと考え、定数として取扱う。右辺第2項の事前確率は、パラメータの分布を正規分布と考え、

$$P(X | \omega_k) = (2\pi)^{-d/2} |W_k|^{-1/2}$$

$$\cdot \exp \{ -1/2 (X - \mu_k)^t \cdot W_k^{-1} \cdot (X - \mu_k) \}$$

【0079】で表わされる。

【0080】

【外15】

$$P(X)$$

【0081】は単語とその周辺情報も含めて、生起し得る全ての入力条件に対する確率の和であり、パラメータがLPCケプストラム係数やバンドパスフィルタ出力の場合は、正規分布に近い分布形状になると考えることができる。

【0082】

【外16】

$$P(X)$$

$$P(X) = (2\pi)^{-d/2} |W_e|^{-1/2}$$

$$\cdot \exp \{ -1/2 (X - \mu_e)^t \cdot W_e^{-1} \cdot (X - \mu_e) \}$$

【0089】(数10)、(数11)を(数9)に代入し、対数をとって、定数項を省略し、さらに-2倍すると、次式を得る。

$$L_k = (X - \mu_k)^t \cdot W_k^{-1} \cdot (X - \mu_k) - (X - \mu_e)^t \cdot W_e^{-1} \cdot$$

$$(X - \mu_e) + \log |W_k| / |W_e|$$

【0091】この(数12)は、ベイズ距離を事後確率化した式であり、識別能力は高いが計算量が多いという欠点がある。この式を次のようにして線形判別式に展開する。全ての単語に対する全ての部分パターンそして周囲パターンも含めて共分散行列が等しいものと仮定する。このような仮定のもとに共分散行列を(数6)によって共通化し、(数12)の

\*【0072】

【外13】

x

【0073】とする（簡単のため当分の間i, jを除いて記述する）。単語kの部分パターン $\omega_k$ に対する事後確率

【0074】

【外14】

$$P(\omega_k | X)$$

【0075】はベイズ定理を用いて次のようになる。

【0076】

【数9】

※【0078】

【数10】

20 ☆【0083】が正規分布に従うと仮定し、平均値を

【0084】

【外17】

$\mu_k$

【0085】、共分散行列を

【0086】

【外18】

$W_e$

【0087】を用いると、(数11)のようになる。

【0088】

30 【数11】

☆

☆【0090】

【数12】

☆

【0092】

【外19】

$W_k^{-1}$

【0093】、

【0094】

【外20】

$W_e$

11

【0095】のかわりに  
 【0096】  
 【外21】

w

$$(X - \mu_k)^t \cdot W_k^{-1} \cdot (X - \mu_k) = X^t \cdot W^{-1} \cdot X - a_k^t \cdot X + b_k$$

【0099】

$$(X - \mu_e)^t \cdot W_e^{-1} \cdot (X - \mu_e) = X^t \cdot W^{-1} \cdot X - a_e^t \cdot X + b_e$$

【0100】(数13)、(数14)において

【0101】

【数15】

$$a_k = 2W^{-1} \cdot \mu_k, b_k = \mu_k^t \cdot W^{-1} \cdot \mu_k$$

【0102】

【数16】

$$a_e = 2W^{-1} \cdot \mu_e, b_e = \mu_e^t \cdot W^{-1} \cdot \mu_e$$

【0103】である。また、(数12)の第3項は0になる。従って、(数12)は次のように簡単な一次判別式になる。

【0104】

【数17】

$$L_k = (b_k - b_e) - (a_k - a_e)^t \cdot X$$

【0105】ここで、改めて、入力第jフレーム成分(数2)と単語kの第iフレーム成分の部分パターンとの距離として(数17)を書き直すと、

【0106】

【数18】

$$L_k^{i,j} = B_k^i - A_k^i \cdot X_j$$

【0107】ここで

【0108】

【外22】

$$g(i, j) = \min \begin{cases} g(i-1, j-2) + \ell(i, j) \\ g(i-1, j-1) + \ell(i, j) \\ g(i-2, j-1) + \ell(i-1, j) + \ell(i, j) \end{cases}$$

【0116】となる。経路判定部6は、(数19)における3つに経路のうち累積距離が最小になる経路を選択する。

【0117】図3は、DP法によって累積距離を求める方法を図示したものである。図のようにペン型非対称のパスを用いているが、その他にもいろいろなパスが考えられる。DP法の他に線形伸縮法を用いることもできるし、また隠れマルコフモデルの手法(HMM法)を用いてもよい。

【0118】このようにして、逐次、距離を累積していく、 $i = I_k$ 、 $j = J$ となる時点での累積距離 $G_k(I$

12

\*【0097】を代入すると、(数12)の第1項、第2項は次のように展開できる。

【0098】

【数13】

【数14】

10※

 $A_k^i$ 

【0109】は(数7)で、

【0110】

【外23】

 $B_k^i$ 

【0111】は(数8)で与えられる。 $L_k^{i,j}$ は、単語kの第i部分パターンと入力第jフレーム近隣のベクトルの部分類似度である。

【0112】図1において距離累積部7は、各単語に対する部分距離を $i = 1 \sim I_k$ の区間に対して累積し、単語全体に対する距離を求める部分である。その場合、入力音声長(Jフレーム)を各単語の標準時間長 $I_k$ に伸縮しながら累積する必要がある。この計算はダイナミックプログラミングの手法(DP法)を用いて効率よく計算できる。

【0113】いま、例えば「サン」の累積距離を求めることにすると、常に $k = 3$ なのでkを省略して計算式を説明する。

【0114】入力第jフレーム部分と第i番目の部分パターンとの部分距離 $L^{i,j}$ を $l(i, j)$ と表現し、

30  $(i, j)$ フレームまでの累積距離を $g(i, j)$ と表現することになると、

【0115】

【数19】

※

40  $k, J$ )を単語ごとに求める。

【0119】判定部8は、累積距離 $G_k(I_k, J)$ の最小値を求めて、(数20)により認識結果

【0120】

【外24】

 $\hat{k}$ 

【0121】を出力する。

【0122】

【数20】

$$\hat{k} = \underset{k}{\operatorname{argmin}} G_k(I_k, J)$$

13

【0123】（実施の形態2）次に、図4に本発明の実施の形態2の音声認識装置の機能ブロック図を示し、説明する。実施形態1では、音声区間検出の後にパターンマッチングを行なったが、実施の形態2では音声区間検出が不要である。入力信号の中から距離が最小の部分を出出すことによって単語を認識する方法であり、「ワードスポッティング法」の1つである。

【0124】この方法は「入力信号中に目的の音声が含まれていれば、その音声の区間において正しい標準パターンとの距離（累積距離）が最小になる」という考え方に基づく方法である。したがって、入力音声の前後のノイズ区間を含む十分長い入力区間において1フレームずつシフトしながら、標準パターンとの照合を行なっていく方法を採用。図4において、図1と同一番号のブロックは同じ機能を持つ。図4が図1と異なる部分は、音声区間検出部9を有しないことと、判定部8のかわりに距離比較部12と一時記憶11が存在することである。以下実施の形態1と異なる部分のみを説明する。

【0125】まず、パターンマッチングが始る時点（ $j=1$ の時点）が音声の始端よりも前にあり、パターンマッチングが終了する時点（ $j=J$ の時点）が音声の終端よりも後にある。パターンマッチングの終了を検出する方法はいろいろと考えられるが、本実施例では全ての標準パターンとの距離が十分大きくなる時点をも  $j=J$  としている。

【0126】標準パターンの作成法は、実施の形態1と全く同じである。ただ、音声サンプルを用いて周囲パターンを作成する範囲は音声区間の前後の十分広い区間を用いる必要がある。その理由は、（数9）の分母項

【0127】

【外25】

 $P(X)$ 

【0128】は、「パターンマッチングの対象となる全てのパラメータに対する確率密度である」という定義によるものである。

【0129】実施の形態1との一番大きな構成上の違いは、単語ごとの累積距離の大小比較をフレームごとに行なう点である。距離比較部12は（数21）により、入力の第  $j$  フレームにおける各単語の累積距離  $G_k(I_k, j)$  を比較して、第  $j$  フレームにおいて累積距離が最小となる単語

【0130】

【外26】

 $\hat{k}_j$ 

【0131】を求める。そして、そのときの最小値も同時に求めておく。即ち、

【0132】

【数21】

14

$$\hat{k}_j = \underset{k}{\operatorname{argmin}} G_k(I_k, j)$$

【0133】

【数22】

$$\hat{G}_j = \underset{k}{\min} G_k(I_k, j)$$

【0134】一時記憶11には  $j-1$  フレームまでに出現した累積距離の最小値  $G_{min}$  と累積距離が最小となった時の標準パターン名  $k$  が記憶されている。

【0135】  $G_{min}$  と

【0136】

【外27】

 $\hat{G}_j$ 

【0137】を比較し、

【0138】

【外28】

$$G_{min} < \hat{G}_j$$

【0139】ならば一時記憶11はそのままにして、次のフレーム（ $j=j+1$ ）へ進む。

【0140】

【外29】

$$G_{min} > \hat{G}_j$$

【0141】ならば、

【0142】

【外30】

$$G_{min} = \hat{G}_j, \hat{k} = \hat{k}_j$$

30

【0143】として次のフレームへ進む。このように、一時記憶11には常にそのフレームまでの最小値と認識結果が残っていることになる。パターンマッチング範囲の終端（ $j=J$ ）に達した時、一時記憶11に記憶されている

【0144】

【外31】

 $\hat{k}$ 

【0145】が認識結果である。実施の形態2は、騒音中の発声など、音声区間検出が難しい場合には有効な方法である。

【0146】本発明の効果を確認するため、男女計150名が発声した10数字データを用いて認識実験を行なった。このうち100名（男女各50名）のデータを用いて標準パターンを作成し、残りの50名を評価した。

評価条件を（表1）に示し、

【0147】

【表1】

15

16

部分パターン作成条件	$L_1=L_2=5, m=1$ (〈数2〉参照)
フレーム周期	10ms
標本化周波数	10kHz
入力フィルタの帯域	300~5000Hz
実施の形態2における評価範囲	音声の始端の前15フレーム~音声の終端後15フレーム

【0148】評価結果を(表2)に示す。

\*【表2】

【0149】

\*

実施の形態1	99.6%(誤りは2つ) 「サン」→「ナナ」, 「ロク」→「ゴ」
実施の形態2	99.4%(誤りは3つ) 3つとも「ロク」→「ゴ」
従来例の方法	97.5%

【0150】このように本実施例における認識率向上は非常に顕著である。

【0151】

【発明の効果】本発明は、複数のフレームで形成される入力ベクトルと、単語音声の部分(標準)パターンとの部分距離を事後確率に基づく統計的距離尺度で求め、フレームをシフトしながら入力ベクトルを更新して各部分ベクトルとの間の距離を累積していき、累積距離を最小とする単語を認識結果とする方法に関するもので、語彙数の増加や騒音に対して頑強で高い認識率が得られるという効果が得られる。

【0152】さらに、計算の方法が単純であるので信号処理プロセッサ(DSP)を用いた小型装置として容易に実現できる。

【0153】このように本発明は実用上有効な方法であり、その効果は大きい。

【図面の簡単な説明】

【図1】本発明の実施の形態1における音声認識装置の機能ブロック図

【図2】本発明における標準パターン作成法における部

分パターン、周囲パターン作成法を説明する概念図

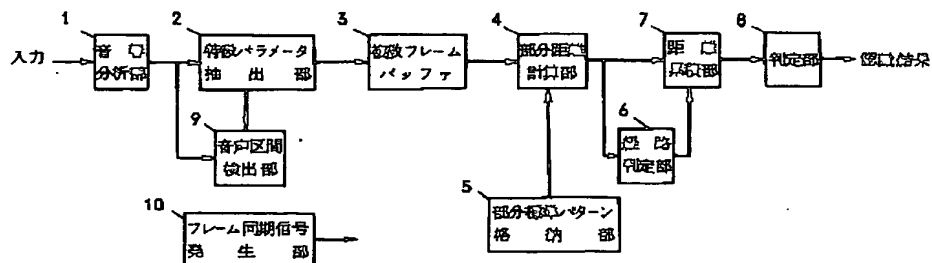
【図3】本発明における入力音声と部分パターンを接続した標準パターンの照合をダイナミックプログラミング法で計算する方法を示した模式図

【図4】本発明の実施の形態2における音声認識装置の機能ブロック図

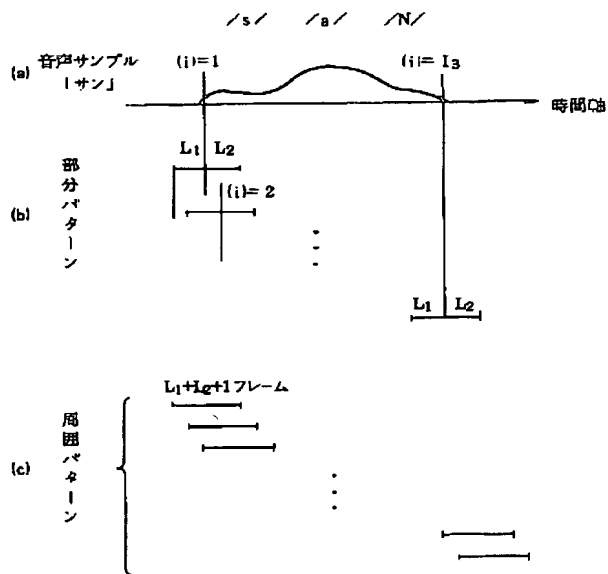
【符号の説明】

- 1 音響分析部
- 2 特徴パラメータ抽出部
- 3 複数フレームバッファ
- 4 部分距離計算部
- 5 部分標準パターン格納部
- 6 経路判定部
- 7 距離累積部
- 8 判定部
- 9 音声区間検出部
- 10 フレーム同期信号発生部
- 11 一時記憶
- 12 距離比較部

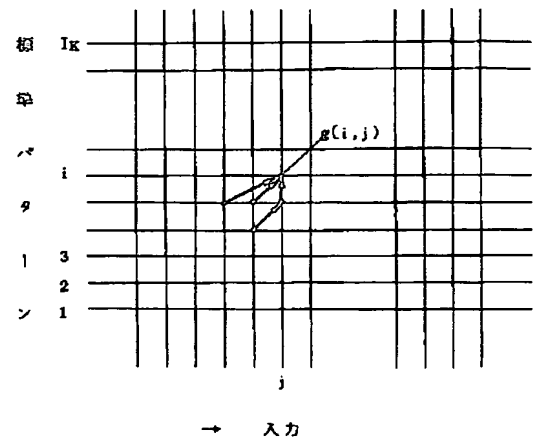
【図1】



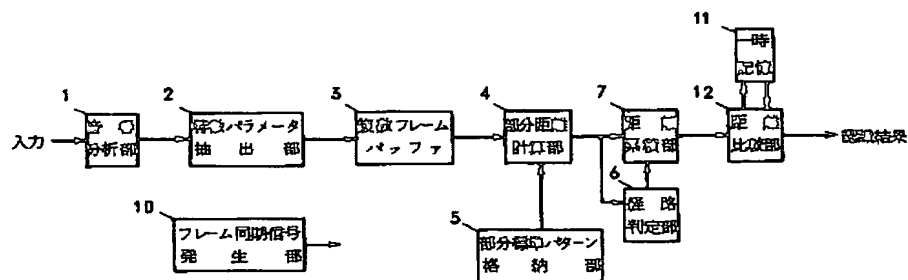
【図 2】



【図 3】



【図 4】



フロントページの続き

(72) 発明者 木村 達也  
 神奈川県川崎市多摩区東三田 3 丁目 10 番 1  
 号 松下技研株式会社内